

Genetic Algorithm Wrappers for Feature Subset Selection in Supervised Inductive Learning

William H. Hsu

Cecil P. Schmidt

James A. Louis

Laboratory for Knowledge Discovery in Databases, Kansas State University
234 Nichols Hall, Manhattan, KS 66506

{bhsu | cps4444 | jal8334}@cis.ksu.edu

<http://www.kddresearch.org>

We derive a validation-based genetic algorithm for feature selection in supervised inductive learning, based upon the following loss functions:

1. **Inferential loss:** Quality of the model produced by an inducer as detected through inferential loss evaluated over a holdout validation data set $D_{val} \equiv D \setminus D_{train}$
2. **Model loss:** “Size” of the model under a specified coding or representation
3. **Ordering loss:** Inference/classification-independent and model-independent measure of data quality given only training and validation data D and hyperparameters \mathbf{a}

$$f(\mathbf{a}, D, \bar{\mathbf{I}}_e) = a \cdot f_a(\mathbf{a}, D, \bar{\mathbf{I}}_e) + b \cdot f_b(\mathbf{a}, D) + c \cdot f_c(\mathbf{a}, D) \quad (1)$$

$$f_a^{BN}(\mathbf{a}, D, \bar{\mathbf{I}}_e) = 1 - \sqrt{\frac{1}{\sum_{X_i \in \mathbf{X} \setminus E} a_i \sum_{X_j \in \mathbf{X} \setminus E} \sum_{j=1}^{a_i} (P'(x_{ij}) - P(x_{ij}))^2}} \quad (2)$$

$$f_a^{DT}(\mathbf{a}, D) = 1 - \frac{m_{correct}}{m_{val}} \quad (3)$$

where $m_{correct} \equiv h.classification-accuracy(D_{val}, select(\mathbf{a}))$

$h \equiv h_0.train(D_{train}, select(\mathbf{a}))$

$m_{val} \equiv |D_{val}|$

$$f_b^{BN}(\mathbf{a}, D) = 1 - \frac{\sum_{i=1}^n (a_i \cdot \max_{|X_j \in Pa_{s_i}|} \prod a_j, 1)}{\prod_{i=1}^n a_i} \quad (4)$$

where $a_i \equiv arity(X_i, B = (\mathcal{X}, E, \Theta))$

$(E, \Theta) = K2(\mathbf{a}, D_{train})$

$$f_b^{DT}(\mathbf{a}, D) = 1 - \frac{h.size(\cdot)}{s_{max}} \quad \text{e.g., } s_{max} = m \quad (5)$$

$$f_c^{DT}(\mathbf{a}) = 1 - \frac{|\mathbf{a}|}{n} \quad (6)$$

$$a + b + c = 1 \quad (7)$$

In related work on genetic wrappers for variable selection in supervised inductive learning, we adapted Equation (3) [HWRC02] from similar fitness functions developed by Cherkauer and Shavlik for decision tree pre-pruning and by Guerra-Salcedo and Whitley for connectionist learning

[GW99]. This breadth of applicability demonstrates the generality of simple genetic algorithms as wrappers for performance tuning in supervised inductive learning.

In experiments using the UC Irvine Machine Learning Database repository, this system is shown to be competitive with search-based feature selection wrappers.

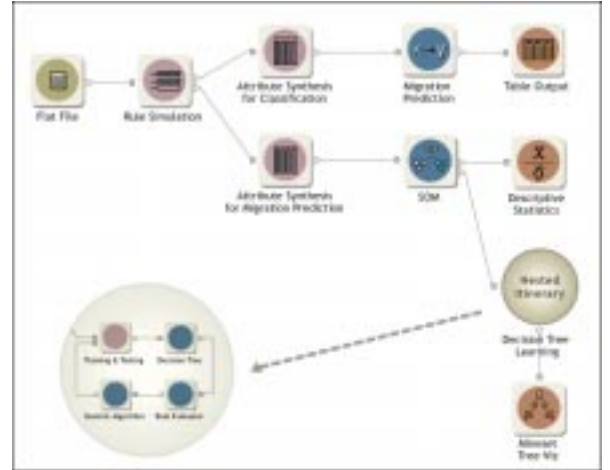


Figure 1. Itinerary for MLJ-CHC

Figure 1 illustrates a real-world application [HWRC02] – automobile insurance risk analysis – that uses the GA wrapper system (depicted in the lower-left inset). Preliminary results on this test bed also indicate that the system is competitive with search-based wrappers.

References

- [Be90] D. P. Benjamin, editor. *Change of Representation and Inductive Bias*. Kluwer Academic Publishers, Boston, 1990.
- [GW99] C. Guerra-Salcedo and D. Whitley. Genetic Approach to Feature Selection for Ensemble Creation. In *Proceedings of the 1999 International Conference on Genetic and Evolutionary Computation (GECCO-99)*. Morgan-Kaufmann, San Mateo, CA, 1999.
- [HWRC02] W. H. Hsu, M. Welge, T. Redman, and D. Clutter. Constructive Induction Wrappers in High-Performance Commercial Data Mining and Decision Support Systems. *Knowledge Discovery and Data Mining*, Kluwer, 2002.