

Feature Selection For Self-Supervised Learning*

Pierre Dangauthier and Pierre Bessière and Anne Spalanzani

E-Motion

INRIA / GRAVIR - CNRS

655 Avenue de l'Europe, Montbonnot

38334 Saint Ismier cedex - France

pierre.dangauthier@imag.fr

Abstract

A foundation of the developmental approach to robotics is that learning must be grounded on sensorimotor interaction. In order to behave autonomously, a robot has to build its own model of the world by searching and exploiting statistical regularities in his sensorimotor domain. Self-supervised learning consists in relying on previous knowledge to acquire new skills. We propose to mix self-supervised learning with our probabilistic programming method, the Bayesian Robot Programming Framework. This idea corresponds to achieve feature selection for searching for relevant sensors. We compare several feature selection algorithms and validate them on a real robotic experiment.

Introduction

In a real environment, a robot needs to continuously improve its knowledge to interact with humans. It needs to better the performance of its previously known skills and to learn new ones. In a complex environment, learning totally new behaviors probably require human feedback. But there are simple situations where unsupervised, or more precisely self-supervised learning, is possible. We study in this article the case of a robot which improves the use of its body without any human supervision. In a simple tracking task, the robot learns by itself to use its laser sensor (SICK) instead of its camera. This is done by searching for correlations in the space of its sensor and motor variables during the tracking behavior.

Self-supervised learning

Self-supervised learning is biologically inspired and is strongly related to mental development (Weng *et al.* 2001). This way of learning takes place during baby's development, but also during adulthood. For instance, a beginner strongly relies on his sight to use a computer keyboard. Then, he

learns to use his finger sensibility and spatial position to perform the same task. One can say that he learned a new sensorimotor behavior thanks to the supervision of his sight.

This learning is possible by searching, and exploiting statistical regularities in the sensorimotor space. The exploitation of statistical regularities can be seen as the foundation of learning and is a promising model of cognition (Barlow 2001). Finding regularities means finding compact data representations, with which a system becomes able to generalize and makes sense of its perceptions.

Bayesian Robot Programming

Apart from the notion of auto-supervised learning, a domestic robot perceives abundant, uncertain and often contradictory information. Sensors and actuators are not perfectly reliable. Therefore, classical determinist programming has been shown to be unable to address real world problems (Bessière *et al.* 1998). We prefer to use a probabilistic approach method called **Bayesian Programming**. Bayesian Robot Programming is a promising candidate in this context (Bessière & the LAPLACE research Group 2003) and has given several interesting results (Lebeltel *et al.* 2003). Moreover, Bayesian Inference is a promising model for understanding animal perception and cognition (BIBA 2001 2005).

Bayesian Robot Programming is based on the subjective interpretation of probabilities deeply described by E.T. Jaynes (Jaynes 2003). A Bayesian program is a probabilistic representation of the relations between sensors and actuators in order to perform a specified task.

In this framework, a Bayesian programmer starts to describe a task by specifying **preliminary knowledge** to the robot. Then the robot processes Bayesian inference in order to take a decision regarding its inputs.

The reader can find detailed examples of non-trivial Bayesian programs in (Lebeltel *et al.* 2003).

Goal

In this work, we present an attempt to automate the creation of a Bayesian program. The "relevant variables" part of a new program will be autonomously discovered under the supervision of another Bayesian program.

The initial behavior of the robot is to track a red ball with a digital camera. This behavior is controlled by an initial

*This work is partially supported by the European Bayesian Inspired Brain and Artefacts (BIBA) project, by the French National Institute for Research in Computer Science and Control (INRIA), and by the French Research Ministry.
Copyright © 2006, American Association for Artificial Intelligence (www.aaai.org). All rights reserved.

Bayesian program. During a tracking experiment, all sensor inputs are recorded, including the direction of the ball. Then, the robot looks for sensors highly correlated with the position of the ball. After that, the robot builds a new tracking program with these sensors, and is then able to stare at the target without its camera.

We present here several algorithms for discovering the relevant variables related to that simple tracking task.

Approach

Once the robot has recorded, for each time step, the position of the ball given by its camera and all other sensor values, the problem is reduced to classical supervised learning. Thus, the goal is to find the minimal set of variables allowing a good prediction of the position.

Finding relevant variables for good predictions is a well known problem in the AI field. In our case, the selected subset of variables will be the input of a naive Bayesian classifier (Domingos & Pazzani 1997). This search is therefore called **feature selection for machine learning**.

This article presents different feature selection methods for finding correlations between a variable and the sensor data. We show that our methods enhance the computing time and the recognition rate of the classification. A study of the selected sensors allows us to validate the relevance of our approach regarding the investigation for new sensorimotor modalities.

Experiment

Presentation

The problem is as follows: we have a robot with a lot of different sensors (laser range SICK, proximeters, thermometer, odometers, battery level...) and a ball is moving in the horizontal plane. The ball position is defined by the angle θ between the ball and the robot's main axe. A set of learning examples is recorded. For each time step, we record θ and the related values of all the sensors. After a learning stage, our robot should be able to guess the position of the ball knowing its sensor values.

The new constructed Bayesian program will be:

- *Relevant variables* are sensors carrying information about θ . We denote sensors variables by $X_i, \forall i \in [0 \dots N - 1]$.
- *Decomposition*: The decomposition of the joint probability distribution we have chosen is called **naive Bayes model**. In our framework, it is reasonable to assume that sensor values are conditionally independent given the value of θ . Thus the decomposition is:

$$P(X_0 \dots X_{N-1}, \theta) = P(\theta) \prod_{i=0}^{N-1} P(X_i|\theta).$$

- *Priors*: We choose Laplace's histograms (Jaynes 2003). Laplace's histograms, which are simple histograms with non-zero values, are frequently used to model discrete probability distributions.
- *Question*: Once this program is defined, the goal for the robot is to infer the position of the ball θ :

$$\begin{aligned} \theta &= \text{Argmax } P(\theta|X_0 \dots X_N) \\ &= \text{Argmax } \frac{P(X_0 \dots X_N|\theta)}{P(X_0 \dots X_N)} \\ &= \text{Argmax } \prod_{i=0}^N P(X_i|\theta). \end{aligned}$$

Robotics

We carried out our experiments on the **BibaBot**, the robot of the **BIBA European project** (BIBA 2001 2005). It is a middle sized (1 - 0.5 - 0.5 m) wheeled robot equipped with a lot of different sensors. In this work, the robot stay motionless.

It is equipped with:

- A Pan-Tilt camera,
- A laser scanner (SICK LMS 200) that scans its surroundings in 2D. It gives 361 measures of distance in the horizontal plane. We will consider each laser beam as an individual sensor.
- 3 odometry sensors,
- 4 bumpers,
- 8 ultrasonic sensors,
- 15 infrared proximity sensors.
- a clock and the battery voltage.

The experiment is then done in three steps:

- Firstly the visual tracking program is launched and a red ball is moving in front of the robot. The program makes the Pan-Tilt camera follow the ball. Thus, we can record θ as the angle of the Pan axis. In the same time we record all the other sensor values.
- Secondly, our feature selection algorithms are launched off line on the collected data. For a given algorithm, a subset of relevant sensors is found.
- Finally we launch the new Bayesian program and see if the robot can guess the position of the ball, or of another object, without any visual information from its camera.

The recognition rate of the position determines the quality of the sensor selection algorithm.

Validation criteria

We have different criteria to judge the pertinence of our sensor selection algorithms. The first one is obviously the recognition rate of the classification. Is the robot able to locate the ball with the generated subset of variables? We have also to take into account the computing cost of feature selection, which is crucial for embedded mobile robotics. Another important criterion is the size of the final subset. We aim at minimizing it, while keeping a good recognition rate. The recognition rate is not a monotonic function of the subset size.

Beyond those criteria, we have also to consider the number and the nature of free parameters in our algorithms. In a context of autonomy, the part of the programmer should remain minimal.

State of the art

Feature selection is an active field in computer science, especially for *data mining*. Feature selection for machine learning consists in finding the “best” subset of features (sensors for us) for a classification algorithm. Selecting a subset of features before classifying has several advantages:

- It reduces the dimensionality of the problem allowing the application of more complex algorithms,
- It leads to a simpler model of data (Occam Razor argument (Blumer *et al.* 1987)),
- It enhances the classifier performance, speed and generality,
- It leads to a more comprehensible model and shows conditional dependencies between variables.

Usually feature selection algorithms remove independent or redundant variables.

General feature selection structure

It is possible to derive a common architecture from most of the feature selection algorithms (see Fig. 1). These algorithms create a subset, evaluate it, and loop until an ending criterion is satisfied (Liu & Motoda 1998). Finally the subset found is validated with the classifier algorithm on real data.

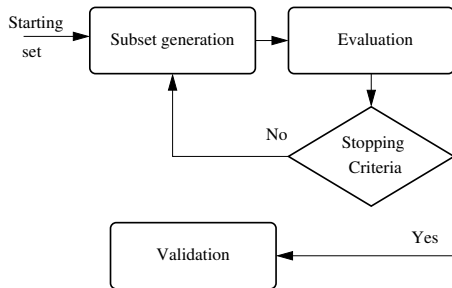


Figure 1: *General feature selection structure*

Subset Generation The subset generation stage is a search procedure in the space of all subsets. As the size of this space is 2^N , exhaustive search methods are often helpless. Non deterministic search like evolutionary search is often used (Yang & Honavar 1998) (Ritthoff *et al.* 2002). It is also possible to employ a heuristic function to build the subsets. There are two main families of heuristic search methods: *forward addition* (Koller & Sahami 1996) (starting with an empty subset, we add features after features by local search) or *backward elimination* (the opposite). Ref.(Koller & Sahami 1996) presents theoretical arguments in favor of backward elimination, but experimental results (John, Kohavi, & Pfleger 1994) shows that forward addition and backward elimination are equivalent. The reader can found in (Liu & Motoda 1998) a detailed nomenclature of feature selection algorithms.

Subset Evaluation

A simple method for evaluating a subset is to consider the performance of the classifier algorithm when it runs with that subset. This way, the classifier algorithm is wrapped in the loop, and the method is classified as a *wrapper*. On the contrary, *filter* methods do not rely on the classifier algorithm, but use other criteria based on correlation notions.

Wrappers Wrappers have been introduced by John, Kohavi and Pfleger in 1994 (John, Kohavi, & Pfleger 1994). Usually, subset evaluations are a compromise between the learning performance and the number of kept features. Thus wrappers generate well suited subsets for recognition tasks because they take into account intrinsic bias of the learning algorithm. Another advantage is their conceptual simplicity. There is no need to really understand causalities in data, they only require generating and testing subsets.

But wrappers have several drawbacks. Firstly they do not clarify conditional dependencies between variables, providing theoretical justification for keeping this or this variable. Since selected subsets are specific to a given classifier algorithm, if it is changed, the selection is not valid anymore. More important, wrappers are computationally costly, so this approach may become intractable.

Filters Filters are quicker and based on theoretical notions. But they often give slightly worse recognition rates. In order to rank a subset, a naive solution consists in giving a mark to each variable independently of others, and to sum those scores. Such a feature ranking method can be done by evaluating correlation between a variable and θ . But (Guyon & Elisseeff 2003) exposes simple examples showing that this is clearly insufficient. For instance, this can not eliminate redundant variables if they are highly correlated with the target.

In contrast an elegant solution is to consider a subset as a whole, as done in the structural learning method for Bayesian networks (Heckerman, Geiger, & Chickering 1994). In our case, this can be reduced to the search of Markov blankets (Koller & Sahami 1996).

But those theoretical methods are NP-complete and only approximations are implemented.

For those reasons, intermediate methods based on statistical analysis (Ghiselli 1964) have proven their efficiency (Hall 1998). The idea of Ghiselli is to give a good mark to a subset if its variables are highly correlated with the target, but slightly correlated between each other. This is summarized in the following formula:

$$r_{\theta S} = \frac{k\overline{r_{\theta i}}}{\sqrt{k + k(k-1)\overline{r_{ij}}}}$$

where $r_{\theta S}$ is the score of the subset, $\overline{r_{\theta i}}$ is the average correlation between the k variable and θ , and $\overline{r_{ij}}$ is the average of the k^2 intra-correlations. This formula is an approximation in the way that we only consider first order dependencies.

Referring to that method, we need a way to evaluate correlation between two variables. While linear correlation used

in (Hall 2000) is clearly inadequate, it is possible to use classical statistics (like χ^2) estimator (Kendall & Stewart 1977)(Liu & Setiono 1995), or an information based measure. Some authors have explored pure Bayesian independence tests (Margartitis & Thrun 2001) (Zaffalon & Hutter 2002), while others rely on other notions as consistency (Almuallim & Dietterich 1994). Recent works try to combine filters and wrappers (Guyon & Elisseeff 2003).

Stopping criteria and validation

Different kinds of stopping criteria can be used: a computing time, a number of kept variables, a heuristic evaluation of the last subset or a recognition rate. In robotics, this criterion should not be a free parameter, for ensuring autonomy and plasticity. The validation of the final subset will be done by calling the learning algorithm.

Algorithms

We have implemented and compared nine different algorithms. Eight of them were recombination of state of the art methods, and one is original. We quickly present here the best four algorithms: a filter and a wrapper based on genetic algorithms (*WrappGA* and *FiltGA*), a filter based on Ghiselli formula (*GhiselliFilt*) and another filter based on the conditional independency hypothesis (*CondIndFilt*).

WrappGA

A genetic algorithm generates subsets which are ranked according to the performance of the classifier algorithms. It ends as soon as a predefined number of generations has been computed.

FiltGA

In this case the subset fitness is computed thanks to Ghiselli formula. Correlation between two variables is estimated in the information theory framework by a cross-entropy estimator. The Kullback-Leibler (KL) (S. & Leibler 1951) distance between two probability distributions defines the mutual information between two variables :

$$I(X_i, X_j) = \sum_{X_i} \sum_{X_j} P(X_i, X_j) \log \left(\frac{P(X_i, X_j)}{P(X_i)P(X_j)} \right).$$

This is in fact the KL distance between $P(X_i, X_j)$ and $P(X_i)P(X_j)$. Therefore the more independent X_i and X_j , the smaller $I(X_i, X_j)$. Taking into consideration that $I(X_i, X_j)$ is biased in favor of variables which have lots of possible values, we can define the uncertain symmetrical coefficient by:

$$USC = 2 \left[\frac{I(X_i, X_j)}{H(X_i) + H(X_j)} \right].$$

Where $H(X)$ is the entropy of X 's distribution. USC is null for independent variables and equals to 1 when X_i and X_j are deterministically linked.

This algorithm ends when a predefined number of generations is computed.

GhiselliFilt

GhiselliFilt is basically the same as *FiltGA*, except that the search method is no more evolutionary, but it is a *forward addition* procedure. We start with an empty set and we add variables one by one. At each step we add the variable which maximizes the score of the candidate subset. The score of a subset is computed with the Ghiselli heuristic, which takes into account correlations with θ and inter correlations.

The algorithm ends when it becomes impossible to add a feature without decreasing the score of the subset.

CondIndFilt

In this method, we explicitly consider that sensor variables are independent knowing the position of the ball. Several robotic experiments have proven the validity of this hypothesis. The search procedure is *backward elimination*. We start with the complete set of variables, and we try to remove the less informative features. In order to choose a variable, we compute the symmetrized KL distance Δ between two probability distributions.

$$\Delta(\mu, \sigma) = D_{KL}(\mu, \sigma) + D_{KL}(\sigma, \mu).$$

If P is the original distribution, and P_i is the distribution considering X_i independent of θ , we have:

$$\begin{aligned} P(X_0 \dots X_{N-1} \theta) &= P(\theta) \prod_{k=0}^{N-1} P(X_k | \theta) \\ P_i(X_0 \dots X_{N-1} \theta) &= P(\theta) P(X_i) \prod_{k=0, k \neq i}^{N-1} P(X_k | \theta). \end{aligned}$$

Then, we just have to compute, for each candidate variable

$$\Delta(P, P_i) = \sum_{\theta} P(\theta) \sum_{X_i} \log \left(\frac{P(X_i)}{P(X_i | \theta)} \right) - (P(X_i) - P(X_i | \theta)).$$

The main drawback is that the algorithm stops when a predefined number of variables are eliminated.

Remark: we have shown that, under conditional independence, the distance relies only on the "difference" between $P(X_i)$ and $P(X_i | \theta)$, e.g. it is not necessary to look at other variables to decide if a variable is useless.

Results

We have tested our algorithms both on simulated and robotic data. Here we present the results of the last four experiments:

- **Exp 1.** The robot learns when a red ball is moved, and the new Bayesian program is tested with the same red ball moved differently.
- **Exp 2.** The robot learns with a red ball, but the new program is tested with a wood cube instead of the ball.

Experiment		WrappGA	FiltGA	GiselliFilt	CondIndFilt	ALL
1	#Sensors	196	118	8	8	392
	Rec.rate	0.143	0.143	0.8	0.42	0,143
	Time	110.5	8.61	3.98	0.72	0,07
2	#Sensors	183	112	8	8	392
	Rec.rate	0.14	0.143	1	0.43	0,1432
	Time	100.4	8.94	3,92	0,726	0,07

Table 2: Results of a simple robotic experiment: This table shows that feature selection hugely improves our Bayesian classifier recognition rate.

Experiment		WrappGA	FiltGA	GiselliFilt	CondIndFilt	ALL
3	#Sensors	14	4	7	8	32
	Rec.rate	0.66	0.2	0.55	0.51	0,55
	Time	2.73	0.15	0.01	0.006	0,002
4	#Sensors	20	5	7	8	32
	Rec.rate	0.40	0.23	0.36	0.38	0,30
	Time	2.97	0.16	0.01	0,005	0,002

Table 3: Same experiments without SICK sensor. One can see that genetic algorithms perform better in this case with a smaller search space.

	WrappGA	FiltGA	GiselliFilt	CondIndFilt
Comp. Time	---	-	++	+++
Reco. Rate	--	--	++	+
#Sensors	--	+	++	parameter
Free Parameters	-	-	++	---

Table 1: Strengths and weaknesses of different algorithms in the general case. "Free Parameters" line represents the number of free parameters, and of the difficulty to choose them.

- **Exp 3.** Same as **Exp 1.**, without the SICK sensor.
- **Exp 4.** Same as **Exp 2.**, without the SICK sensor.

Numerical results for **Exp 1.** and **Exp 2.** are presented in Table 2. We can notice that:

- *GiselliFilt* shows that just 8 laser beams of the SICK are enough to guess the position of the ball. There are two reasons for that:
 - The SICK is a highly reliable sensor.
 - The number of different values chosen during the discretization of θ was quite low (6 different classes for 180 degrees);
- Surprisingly recognition rates are not smaller in Exp.2 than in Exp.1. This shows that the ball and the cube have a similar signature through the kept sensors;
- Algorithms based on genetic algorithms are too slow to be incorporated in a real time processing;
- *WrappGA* and *FiltGA* do not succeed in finding a good subset. The search space is too huge for them; they do not find that only a few sensors are required, even with a lot of generation and a big population. Although in simulated tests, with a dozen of sensors, they find the global optimum better than other methods. Thus we can not use them when too many sensors are used;

- The recognition rate is improved by selecting features. Indeed making a fusion with non relevant sensors decreases the performances.

Numerical results without SICK are presented in Table 3. Removing this sensor involves several consequences:

- As the search space drastically shrinks, genetic search performs better;
- Although the recognition rate among algorithms decreases, a good sensor fusion can find the target quite often;
- Sensors kept in the best subset were mainly I.R. relevant sensors with a few U.S. telemeter.

The relative values of algorithms are presented in table 1. This study helped us to choose *GiselliFilt* as a feature selection algorithm. The *forward addition* search procedure combined with *Giselli* heuristic and an entropic measure of mutual information is a good candidate for real-time feature selection. Indeed, *Giselli* formula is an efficient approximation which highlights sensors highly correlated with the target but slightly correlated with other sensors. Moreover mutual information detects more than only traditional linear correlation.

Conclusion and further work

In this work, we have compared different feature selection algorithms. We have tested some of them to enable a robot to discover autonomously correlations in its sensorimotor domain. In our experiments, the robot found a new way to track a ball, passing from visual to proximity tracking. In this framework, feature selection drastically increases performance and speed of the recognition algorithm. This work is a step toward an automatic search of prior knowledge applied to Bayesian robot programming.

Further work will lead to the integration of the selected algorithm in a permanent self-supervised learning framework. This learning process should be real-time and parallelized.

References

- Almuallim, H., and Dietterich, T. G. 1994. Learning boolean concepts in the presence of many irrelevant features. *Artificial Intelligence* 69(1-2):279–305.
- Barlow, H. 2001. The exploitation of regularities in the environment by the brain. *Behavioral and Brain Sciences* 24(3):602–607.
- Bessi re, P., and the LAPLACE research Group. 2003. Survey: Probabilistic methodology and techniques for artefact conception and development. Technical report, Rapport de recherche INRIA.
- Bessi re, P.; Dedieu, E.; Lebeltel, O.; Mazer, E.; and Mekhnacha, K. 1998. Interpr tation ou description (i) : Proposition pour une th orie probabiliste des syst mes cognitifs sensori-moteurs. *Intellectica* 26-27:257–311.
- BIBA. 2001-2005. Bayesian Inspired Brain and Artefacts European Project. <http://www-biba.inrialpes.fr/>.
- Blumer, A.; Ehrenfeucht, A.; Haussler, D.; and Warmuth, M. K. 1987. Occam’s razor. *Information Processing Letters* 24(6):377–380.
- Domingos, P., and Pazzani, M. J. 1997. On the optimality of the simple bayesian classifier under zero-one loss. *Machine Learning* 29(2-3):103–130.
- Ghiselli, E. E. 1964. *Theory of Psychological Measurement*. McGraw-Hill Book Company.
- Guyon, I., and Elisseeff, A. 2003. An introduction to variable and feature selection. *Journal of Machine Learning Research* 3:1157–1182.
- Hall, M. A. 1998. *Correlation-based Feature Selection for Machine Learning*. Ph.D. Dissertation, Waikato University, Hamilton, NZ.
- Hall, M. A. 2000. Correlation-based feature selection for discrete and numeric class machine learning. In *Proc. 17th International Conf. on Machine Learning*, 359–366. Morgan Kaufmann, San Francisco, CA.
- Heckerman, D.; Geiger, D.; and Chickering, D. M. 1994. Learning bayesian networks: The combination of knowledge and statistical data. In *KDD Workshop*, 85–96.
- Jaynes, E. T. 2003. *Probability Theory : The Logic of Science*. G. Larry Bretthorst.
- John, G. H.; Kohavi, R.; and Pflieger, K. 1994. Irrelevant features and the subset selection problem. In *International Conference on Machine Learning*, 121–129. Journal version in AIJ, available at <http://citeseer.nj.nec.com/13663.html>.
- Kendall, S. M., and Stewart, A. 1977. *The Advanced Theory of Statistics, Volume 1, 4th Edition*. Mcmillan Publishing, New York.
- Koller, D., and Sahami, M. 1996. Toward optimal feature selection. In *International Conference on Machine Learning*, 284–292.
- Lebeltel, O.; Bessi re, P.; Diard, J.; and Mazer, E. 2003. Bayesian robots programming. *Autonomous Robot* 16(1):49–79.
- Liu, H., and Motoda, H. 1998. *Feature Selection for Knowledge Discovery and Data Mining*. Kluwer Academic Publishers.
- Liu, H., and Setiono, R. 1995. Chi2: Feature selection and discretization of numeric attributes. In *Proceedings of 7th IEEE Int’l Conference on Tools with Artificial Intelligence*.
- Margartitis, D., and Thrun, S. 2001. A bayesian multiresolution independence test for continuous variables. In *Uncertainty in Artificial Intelligence: Proceedings of the Seventeenth Conference (UAI-2001)*, 346–353. San Francisco, CA: Morgan Kaufmann Publishers.
- Ritthoff, O.; Klinkenberg, R.; Fischer, S.; and Mierswa, I. 2002. A hybrid approach to feature selection and generation using an evolutionary algorithm. Technical Report CI-127/02, Collaborative Research Center 531, University of Dortmund, Germany.
- S., K., and Leibler. 1951. On information and sufficiency. *Annals of Mathematical Statistics* 22(22):79–86.
- Weng, J.; McClelland, J.; Pentland, A.; Sporns, O.; Stockman, I.; Sur, M.; and Thelen, E. 2001. Autonomous mental development by robots and animals. *Science* 291(291):599–600.
- Yang, J., and Honavar, V. 1998. Feature subset selection using a genetic algorithm. *IEEE Intelligent Systems* 13:44–49.
- Zaffalon, M., and Hutter, M. 2002. Robust feature selection by mutual information distributions. In *Proceedings of the 14th International Conference on Uncertainty in Artificial Intelligence*.