# Learning Qualitative Models from Numerical Data

Jure Žabkar[a,*], Martin Možina[a], Ivan Bratko[a], Janez Demšar[a]

[a] *Artificial Intelligence Laboratory, Faculty of Computer and Information Science, University of Ljubljana, Tržaška 25, 1000 Ljubljana, Slovenia*

## Abstract

Qualitative models are predictive models which describe how changes in values of input variables affect the output variable in qualitative terms, *e.g.* increasing or decreasing. We describe Padé, a new method for qualitative learning which estimates partial derivatives of the target function from training data and uses them to induce qualitative models of the target function. We formulated three methods for computation of derivatives, all based on using linear regression on local neighbourhoods. The methods were empirically tested on artificial and real-world data. We also provide a case study which shows how the developed methods can be used in practice.

*Keywords:* Qualitative Modelling, Regression, Partial Derivatives, Monotone Models

## 1. Introduction

People most often reason qualitatively. For example, playing with a simple pendulum, a five year old child discovers that the period of the pendulum increases if he uses a longer rope. Although most of us are later taught a more accurate numerical model describing the same behaviour, $T = 2\pi\sqrt{l/g}$, we keep relying on the more "operational" qualitative relation in everyday's life. Still, despite Turing's proposition that artificial intelligence should *mimic human intelligence*, not much work has been done so far in trying to learn such models from data.

---

*Corresponding author. Tel.: +386 1 4768 154; fax: +386 1 4768 386.
*Email addresses:* jure.zabkar@fri.uni-lj.si (Jure Žabkar), martin.mozina@fri.uni-lj.si (Martin Možina), ivan.bratko@fri.uni-lj.si (Ivan Bratko), janez.demsar@fri.uni-lj.si (Janez Demšar)

We can formally describe the relation between the period of a pendulum $T$, its length $l$ and the gravitational acceleration $g$ as $T = Q(+l, -g)$, meaning that the period increases with $l$ and decreases with $g$. Different definitions of qualitative relations are discussed in related work. We will base ours on partial derivatives: a function $f$ is in positive (negative) qualitative relation with $x$ over a region $\mathcal{R}$ if the partial derivative of $f$ with respect to $x$ is positive (negative) over the entire $\mathcal{R}$,

$$f = Q_{\mathcal{R}}(+x) \quad \equiv \quad \forall x_0 \in \mathcal{R} : \frac{\partial f}{\partial x}(x_0) > 0 \tag{1}$$

and

$$f = Q_{\mathcal{R}}(-x) \quad \equiv \quad \forall x_0 \in \mathcal{R} : \frac{\partial f}{\partial x}(x_0) < 0. \tag{2}$$

Qualitative models are predictive models which describe qualitative relations between input variables and a continuous output, for instance

$$\text{if } z > 0 \wedge x > 0 \text{ then } f = Q(+x),$$

$$\text{if } z < 0 \vee x < 0 \text{ then } f = Q(-x).$$

For sake of clarity, we omitted specifying the region since it is obvious from the context.

In this paper we propose a new, two-step approach to induction of qualitative models from data. Let the data describe a sampled function given as a set of examples $(\mathbf{x}, y)$, where $\mathbf{x}$ are attributes and $y$ is the function value. In the first step we estimate the partial derivative at each point covered by a learning example. We replace the value of the output $y$ for each example with the sign of the corresponding derivative $q$. Each relabelled example, $(\mathbf{x}, q)$, describes the qualitative behaviour of the function at a single point. In the second step, a general-purpose machine learning algorithm is used to generalize from this relabelled data, resulting in a qualitative model describing the function's (qualitative) behaviour in the entire domain.

Such models describe the relation between the output and a single input variable in dependence of other attributes. For instance, we can model the conditions (*e.g.* public debt, taxation, inflation rate) at which an increase of interest rates will increase/decrease the unemployment. To describe the influence of multiple inputs (*e.g.* the effect of interest rates *and* of inflation and taxation on unemployment) we need to induce multiple models.

The paper includes three major contributions. The first one is the idea of transforming the problem of learning qualitative models to that of learning

ordinary predictive models. Second, we present a new method called Padé[1] for computing partial derivatives from data typical for machine learning. We show three ways of computing the derivatives, all based on variations of local linear regression. Finally, we provide an extensive experimental evaluation of the proposed setup.

## 2. Related Work

The beginnings of the field of qualitative reasoning go back to early work done outside AI. Jeffries and May [1, 2] introduced qualitative stability in ecology, whereas Samuelson [3] discussed the use of qualitative reasoning in economics. The papers by Forbus [4], de Kleer and Brown [5], and Kuipers [6] describe approaches that became the foundations of much of qualitative reasoning work in AI. Kalagnanam *et al.* [7, 8, 9] contributed to the mathematical foundations of qualitative reasoning.

There are a number of approaches to qualitative system identification, also known as learning qualitative models from data. Most of this work is concerned with the learning of QDE (Qualitative Differential Equations) models. One approach is to translate a numerical system's behaviour in time into a qualitative representation, and then check which QDE constraints (or QSIM-type constraints [6]) are satisfied by the qualitative behaviour. The resulting constraints constitute a qualitative model. Examples of this approach are the programs GENMODEL [10, 11], MISQ [12, 13] and QSI [14]. A similar approach can be carried out with a general purpose ILP system (Inductive Logic Programming) to induce a model from qualitative behaviours, which enjoys the advantages of the power and flexibility of ILP. This approach was developed in [15, 16, 17]. Džeroski and Todorovski developed several approaches to discovering dynamics. QMN [18] generates models in the form of qualitative differential equations. LAGRANGE [19] and LAGRAMGE [20] can describe the observed behaviour in the form of differential and partial differential equations respectively. SQUID [21] finds models in terms of semi-quantitative differential equations. Gerçeker and Say [22] fit polynomials to numerical data and use them to induce qualitative models. We believe that their algorithm LYQUID can also be used in static systems although they experiment only with dynamic systems.

---

[1] An acronym for "partial derivative", and the name of a famous French mathematician

Qualitative models of dynamic systems are not necessarily based on QDE. For example, the KARDIO model of the heart [23] uses symbolic descriptions in logic instead. QuMas [24, 25] learned such models from example qualitative behaviours in time and used a kind of ILP for that.

In this paper we only consider static systems. Systems that learn dynamic qualitative models can, to some extent also learn static models by simply not using the functionality for handling temporal dependencies. However, since such systems are intended for dynamic models, their capability in handling static relations are limited, for example to searching for only those static constraints that hold over the whole problem space. The QUIN algorithm [26, 27, 28] is specifically intended for learning qualitative models of static systems, and is therefore the most relevant to the present paper. QUIN induces qualitative trees which are similar to classification trees but have different leaf nodes. The leaves of a qualitative tree contain qualitative constraints (*e.g.* $z = M^{+-}(x, y)$, meaning that $z$ increases with $x$, decreases with $y$ and depends on no other variables). QUIN constructs such trees by computing the qualitative change vectors between all pairs of points in the data and then recursively splitting the space into regions which share common qualitative properties, such as the one given above. Despite similarities, Padé and QUIN are significantly different in that Padé acts as a preprocessor of numerical data and can be used in combination with any attribute-value learning system. Padé computes qualitative partial derivatives in all example data points, and these derivatives become class values for the subsequent learning. For example, Padé combined with a decision tree learner would produce qualitative trees similar to those induced by QUIN. However, Padé combined with a rule learner, say, produces a model in the form of "qualitative rules". The main algorithmic difference between the two systems is that Padé computes quantitative and qualitative partial derivatives in every example point, whereas QUIN computes the degree of consistency of a subregion of numerical data with (essentially) every possible qualitative monotonicity constraint.

Padé differs from all the above methods in being, to our knowledge, the only algorithm for computing (partial) derivatives on point cloud data. Another important difference between Padé and above mentioned algorithms is also that Padé is essentially a preprocessor while other algorithms induce a model. Padé merely augments the learning examples with additional labels, which can later be used by appropriate algorithms for induction of classification or regression models, or for visualization. This results in a number of Padé's advantages. For instance, most other algorithms for learning qualitative models only handle

numerical attributes. Since Padé can be combined with any machine learning method, the learner can also use discrete attributes at no extra cost.

Outside artificial intelligence, there are several numerical methods for computation of derivatives at a given point. Such numerical derivatives could be used for modelling qualitative behaviour. Unfortunately, numerical differentiation computes the function value at points chosen by the algorithm, while we are limited to a given data sample.

## 3. Computation of Partial Derivatives

We will denote a learning example as $(\mathbf{x}, y)$, where $\mathbf{x} = (x_1, x_2, \ldots x_n)$ and $y$ is the value of the unknown sampled function, $y = f(\mathbf{x})$.

We will introduce three methods for estimation of partial derivative of $f$ at point $\mathbf{x_0}$. The simplest one assumes that the function is linear in a small hyper-sphere around $\mathbf{x_0}$ (Figure 1a). It computes *locally weighted linear regression* on examples lying in the hyper-sphere and considers the computed coefficients as partial derivatives. The second method, *$\tau$-regression*, computes a single partial derivative at a time. To avoid the influence of other arguments of the function, it considers only those points in the sphere which lie in a hyper-tube along the axis of differentiation (Figure 1b). The derivative can then be computed with weighted univariate regression. Finally, the *parallel pairs* method replaces the single hyper-tube with a set of pairs aligned with the axis of differentiation (Figure 1c), which allows it to focus on the direction of differentiation without decreasing the number of examples considered in the computation.

### 3.1. Locally weighted regression

Let $\mathcal{N}(\mathbf{x_0})$ be a set of examples $(\mathbf{x_m}, y_m)$ such that $x_{mi} \approx x_{0i}$ for all $i$ (Figure 1a). According to Taylor's theorem, a differentiable function is approximately linear in a neighbourhood of $\mathbf{x_0}$,

$$f(\mathbf{x_m}) = f(\mathbf{x_0}) + \nabla f(\mathbf{x_0}) \cdot (\mathbf{x_m} - \mathbf{x_0}) + R_2. \tag{3}$$

Our task is to find the vector of partial derivatives, $\nabla f(\mathbf{x_0})$. We can solve this as a linear regression problem by rephrasing (3) as a linear model

$$y_m = \beta_0 + \boldsymbol{\beta}^{\mathrm{T}}(\mathbf{x_m} - \mathbf{x_0}) + \epsilon_m, \quad (\mathbf{x_m}, y_m) \in \mathcal{N}(\mathbf{x_0}), \tag{4}$$

where the task is to find $\boldsymbol{\beta}$ (and $\beta_0$) with the minimal sum of squared errors $\epsilon_m$ over $\mathcal{N}(\mathbf{x_0})$. The error term $\epsilon_m$ covers the remainder of the Taylor expansion, $R_2$, as well as noise in the data.
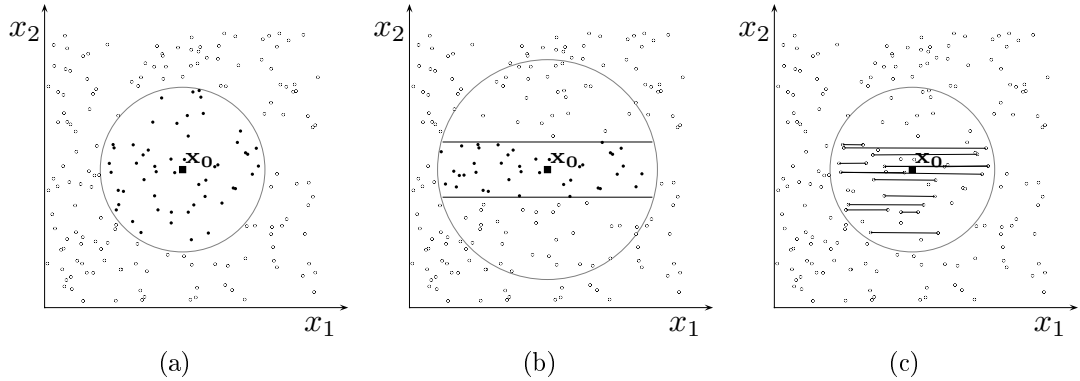
Figure 1: The neighbourhoods for locally weighted regression, $\tau$-regression and parallel pairs.

The size of the neighbourhood $\mathcal{N}(\mathbf{x_0})$ should reflect the density of examples and the amplitude of noise. Instead of setting a predefined radius (*e.g.* $||\mathbf{x_m} - \mathbf{x_0}|| < \delta$), we consider a neighbourhood of $k$ nearest examples and weigh the points according to their distance from $\mathbf{x_0}$,

$$w_m = e^{-||\mathbf{x_m}-\mathbf{x_0}||^2/\sigma^2}. \tag{5}$$

The parameter $\sigma^2$ is fitted so that the farthest example has a negligible weight of $0.001$.

This transforms the problem into locally weighted regression (LWR) [29], where the regression coefficients represent partial derivatives,

$$\begin{bmatrix} \beta_0 \\ \boldsymbol{\beta} \end{bmatrix} = (X^{\mathrm{T}}WX)^{-1}X^{\mathrm{T}}WY, \tag{6}$$

where

$$X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{k1} & \dots & x_{kn} \end{bmatrix} \qquad W = \begin{bmatrix} w_1 & 0 & \cdots \\ 0 & \ddots & \\ \vdots & & w_k \end{bmatrix} \qquad Y = \begin{bmatrix} y_1 \\ \vdots \\ y_k \end{bmatrix} \tag{7}$$

The computed $\boldsymbol{\beta}$ estimates the vector of partial derivatives $\nabla f(\mathbf{x_0})$.

As usual in linear regression, the inverse in (6) can be replaced by pseudo-inverse to increase the stability of the method.

6

### 3.2. $\tau$-regression

The $\tau$-regression algorithm differs from LWR in the shape of the neighbourhood of the reference point. It starts with $\kappa$ examples in a hyper-sphere, which is generally larger than that for LWR, but then keeps only $k$ examples that lie closest to the axis of differentiation (Figure 1b). Let us assume without loss of generality that we compute the derivative with regard to the first argument $x_1$. The neighbourhood $\mathcal{N}(\mathbf{x_0})$ thus contains $k$ examples with the smallest distance $||\mathbf{x_m} - \mathbf{x_0}||_{\backslash 1}$ chosen from the $\kappa$ examples with the smallest distance $||\mathbf{x_m} - \mathbf{x_0}||$, where $|| \cdot ||_{\backslash i}$ represents the distance computed over all dimensions except the $i$-th.

With a suitable selection of $\kappa$ and $k$, we can assume $x_{m1} - x_{01} \gg x_{mi} - x_{0i}$ for all $i > 1$ for most examples $\mathbf{x_m}$. If we also assume that partial derivatives with regard to different arguments are comparable in size, we get $\partial f/\partial x_1(x_{m1} - x_{01}) \gg \partial f/\partial x_i(x_{mi} - x_{0i})$ for $i > 1$. We can thus omit all dimensions but the first from the scalar product in (3):

$$f(\mathbf{x_m}) = f(\mathbf{x_0}) + \frac{\partial f}{\partial x_1}(\mathbf{x_0})(x_{m1} - x_{01}) + R_2 \tag{8}$$

for $(\mathbf{x_m}, y_m) \in \mathcal{N}(\mathbf{x_0})$. We again set up a linear model,

$$y_m = \beta_0 + \beta_1(x_{m1} - x_{01}) + \epsilon_m, \tag{9}$$

where $\beta_1$ approximates the derivative $\frac{\partial f}{\partial x_1}(\mathbf{x_0})$. The task is to find the value of $\beta_1$ which minimizes error over $\mathcal{N}(\mathbf{x_0})$.

Examples in $\mathcal{N}(\mathbf{x_0})$ are weighted according to their distance from $\mathbf{x_0}$ along the axis of differentiation,

$$w_m = e^{-(x_{m1} - x_{01})^2/\sigma^2}, \tag{10}$$

where $\sigma$ is again set so that the farthest example has a weight of 0.001.

The described linear model can be solved by weighted univariate linear regression over the neighbourhood $\mathcal{N}(\mathbf{x_0})$,

$$\frac{\partial f}{\partial x_1}(\mathbf{x_0}) = \beta_1 = \frac{\sum_{x_m \in \mathcal{N}(\mathbf{x_0})} w_m x_{m1} y_m}{\sum_{x_m \in \mathcal{N}(\mathbf{x_0})} w_m x_{m1}^2}. \tag{11}$$

### 3.3. Parallel pairs

Consider two examples $(\mathbf{x_m}, y_m)$ and $(\mathbf{x_l}, y_l)$ which are close to $\mathbf{x_0}$ and aligned with the $x_1$ axis, $||x_m - x_l|| \approx |x_{m1} - x_{l1}|$. Under these premises we can suppose that both examples correspond to the same linear model (9) with the same coefficients $\beta_0$ and $\beta_1$. Subtracting (9) for $y_m$ and $y_l$ gives

$$y_m - y_l = \beta_1(x_{m1} - x_{l1}) + (\epsilon_m - \epsilon_l) \tag{12}$$

or

$$y_{(m,l)} = \beta_1 x_{(m,l)1} + \epsilon_{(m,l)}, \tag{13}$$

where $y_{(m,l)} = y_m - y_l$ and $x_{(m,l)1} = x_{m1} - y_{l1}$. The difference $y_m - y_l$ is therefore linear with the difference in the attribute values $x_{m1}$ and $x_{l1}$

Coefficient $\beta_1$ again approximates the derivative $\frac{\partial f}{\partial x_1}(\mathbf{x_0})$. Note that the model has no intercept term, $\beta_0$.

To compute the derivative using (12) we take the spherical neighbourhood like the one from the first method, LWR. For each pair we compute its alignment with the $x_1$ axis using a scalar product with the base vector $\mathbf{e_1}$,

$$\alpha_{(m,l)} = \left| \frac{(\mathbf{x_m} - \mathbf{x_l})^{\mathrm{T}} \mathbf{e_1}}{||\mathbf{x_m} - \mathbf{x_l}|| \, ||\mathbf{e_1}||} \right| = \frac{|x_{m1} - x_{l1}|}{||\mathbf{x_m} - \mathbf{x_l}||} \tag{14}$$

We select the $k$ best aligned pairs from $\kappa$ points in the hyper-sphere around $\mathbf{x_0}$ (Figure 1c) and assign them weights corresponding to the alignment,

$$w_{(m,l)} = e^{-\alpha_{(m,l)}^2/\sigma^2}, \tag{15}$$

with $\sigma^2$ set so that the smallest weight equals 0.001.

The derivative is again computed using univariate linear regression,

$$\frac{\partial f}{\partial x_1}(\mathbf{x_0}) = \beta_1 = \frac{\sum_{(x_m,x_l) \in \mathcal{N}(\mathbf{x_0})} w_{(m,l)}(x_{m1} - x_{l1})(y_m - y_l)}{\sum_{(x_m,x_l) \in \mathcal{N}(\mathbf{x_0})} w_{(m,l)}(x_{m1} - x_{l1})^2} \tag{16}$$

## 4. Experiments

We first performed an extensive set of experiments on artificially constructed problems. These data sets are used to test Padé with respect to accuracy, the effect of irrelevant attributes and the effect of noise in the data.

To assess the accuracy of induced models, we compute the derivatives and the model from the entire data set. We then check whether the predictions of

8

the model match the analytically computed partial derivatives. We define the accuracy of Padé as the proportion of examples with correctly predicted qualitative partial derivatives. Note that this procedure does not require separate training and testing data set since the correct answer with which the prediction is compared is not used in induction of the model. Where not stated otherwise, experimental results represent averages of ten trials.

Another set of experiments was performed on real-world data. Since the ground truth (the correct partial derivative) on such data sets is not necessarily known, we can not compute the predictive accuracy of the model. In such cases, stability is often used as a measure of quality of a machine learning method [30]. Stability measures the influence that variation of the input has on the output of a system. The motivation for such analysis is to design robust systems that will not be affected by noise and randomness in sampling.

In experiments with the parallel pairs method we reduced the number of arguments to be fitted by setting $\kappa = k$. Since the first batch of experiments on artificial data sets suggested that parameter settings do not significantly affect the performance of methods, we use constant values in other experiments: we use $k = 20$ for LWR and parallel pairs, and $\kappa = 100$ and $k = 20$ for $\tau$-regression, except for real-world data sets where we decreased $\kappa$ to 50 due to a smaller number of examples.

For LWR, we used ridge regression to compute the ill-posed inverse in (6). We used C4.5 as reimplemented in Orange [31] for induction of qualitative models. We preferred this algorithm over potentially more accurate methods, like SVM, since one of our goals was to induce comprehensible models. Besides that, most artificial data sets are simple enough to require only a single-node tree or a tree with two leaves only. C4.5 was run with default settings except where noted otherwise.

## 4.1. Accuracy on artificial data sets

We performed experiments with three mathematical functions: $f(x, y) = x^2 - y^2$, $f(x, y) = x^3 - y$, $f(x, y) = \sin x \sin y$. We sampled them uniform randomly in 1000 points in the range $[-10, 10] \times [-10, 10]$.

Function $f(x, y) = x^2 - y^2$ is a standard test function often used in [27]. Its partial derivative w.r.t. $x$ is $\partial f / \partial x = 2x$, so $f = Q(+x)$ if $x > 0$ and $f = Q(-x)$ if $x < 0$. Since the function's behaviour with respect to $y$ is similar, we observed only results for $\partial f / \partial x$. The accuracy of all methods is close to 100% (Table 1). Changing the values of parameters has no major effect except for $\tau$ regression, where short ($\kappa = 30$ and $\kappa = 50$) and wide ($k \geq 10$) tubes give better accuracy

9

|         | LWR  | Pairs | $\tau$-regression | | | |
|         |      |       | $\kappa = 30$ | 50 | 70 | 100 |
|---------|------|-------|-------|------|------|------|
| $k = 10$ | .991 | .986 | .948 | .968 | .971 | .980 |
| 20      | .993 | .992 | .929 | .960 | .972 | .981 |
| 30      | .992 | .992 | .909 | .953 | .969 | .978 |
| 40      | .993 | .993 |      | .950 | .964 | .974 |
| 50      | .994 | .993 |      | .935 | .961 | .972 |

(a) Accuracy of qualitative derivatives.

|         | LWR  | Pairs | $\tau$-regression | | | |
|         |      |       | $\kappa = 30$ | 50 | 70 | 100 |
|---------|------|-------|-------|------|------|------|
| $k = 10$ | .997 | .993 | .990 | .993 | .990 | .991 |
| 20      | .996 | .995 | .983 | .991 | .986 | .991 |
| 30      | .996 | .994 | .983 | .987 | .988 | .993 |
| 40      | .995 | .995 |      | .978 | .978 | .990 |
| 50      | .998 | .995 |      | .977 | .987 | .991 |

(b) Accuracy of qualitative models induced by C4.5.

Table 1: Results of experiments with $f(x, y) = x^2 - y^2$.

while for very long tubes ($\kappa = 100$) the accuracy decreases with $k$. The latter can indicate that longer tubes reach across the boundary between the positive and negative values of $x$. Induced tree models have the same high accuracy.

Function $f(x, y) = x^3 - y$ is globally monotone, increasing in $x$ and decreasing in $y$ in the whole region. The function is interesting because its much stronger dependency on $x$ can obscure the role of $y$. All methods have a 100% accuracy with regard to $x$. Prediction of function's behaviour w.r.t. $y$ proves to be difficult: the accuracy of $\tau$-regression is 50–60% and the accuracy of LWR is just over 70% (Table 2). Parallel pairs seem undisturbed by the influence of $x$ and estimate the sign of $\partial f / \partial y$ with accuracy of more than 95% with proper parameter settings.

An interesting observation here is that the accuracy of induced qualitative tree models highly exceeds that of point-wise partial derivatives. For instance, qualitative models for derivatives by LWR reach 95–100% despite the low, 70% accuracy of estimates of the derivative. When generalizing from labels denoting qualitative derivatives, incorrect labels are scattered randomly enough that C4.5 recognizes them as noise and induces a tree with a single node.

| | LWR | Pairs | τ-regression κ = 30 | 50 | 70 | 100 |
|---|---|---|---|---|---|---|
| k = 10 | .727 | .690 | .597 | .626 | .639 | .647 |
| 20 | .725 | .792 | .576 | .593 | .619 | .646 |
| 30 | .751 | .879 | .545 | .571 | .609 | .618 |
| 40 | .740 | .919 | | .556 | .588 | .613 |
| 50 | .725 | .966 | | .541 | .558 | .612 |

(a) Accuracy of qualitative derivatives.

| | LWR | Pairs | τ-regression κ = 30 | 50 | 70 | 100 |
|---|---|---|---|---|---|---|
| k = 10 | 1.00 | .898 | .737 | .880 | .890 | .864 |
| 20 | 1.00 | .930 | .745 | .743 | .811 | .860 |
| 30 | .978 | .954 | .752 | .599 | .813 | .777 |
| 40 | .971 | .964 | | .780 | .732 | .700 |
| 50 | .956 | .986 | | .906 | .805 | .760 |

(b) Accuracy of qualitative models induced by C4.5.

Table 2: Results of experiments with $f(x, y) = x^3 - y$.

Function $f(x, y) = \sin x \sin y$ has partial derivatives $\partial f / \partial x = \cos x \sin y$ and $\partial f / \partial y = \cos y \sin x$, which change their signs multiple times in the observed region. The accuracy of all methods is mostly between 80 and 90 percent, degrading with larger neighbourhoods (Table 3). However, the accuracy of C4.5 barely exceeds 50% which we would get by making random predictions. Rather than a limitation of Padé, this shows the (expected) inability of C4.5 to learn this checkboard-like concept.

Since qualitative modelling has been traditionally interested in examples from physics, we also tested Padé's ability to rediscover three simple physical laws from computer generated data. Each domain is described by an equation which was used to obtained 1000 uniform random samples.

*Centripetal force.* The centripetal force $F$ on an object of mass $m$ moving at a speed $v$ along a circular path with radius $r$ is

$$F = \frac{mv^2}{r}.$$

We prepared a data set for this target concept with $m$ sampled randomly from $[0.1, 1]$, $r$ from $[0.1, 2]$ and $v$ from $[1, 10]$.

| | LWR | Pairs | $\tau$-regression $\kappa = 30$ | 50 | 70 | 100 |
|---|---|---|---|---|---|---|
| $k = 10$ | .882 | .858 | .863 | .885 | .890 | .886 |
| 20 | .870 | .841 | .820 | .865 | .880 | .882 |
| 30 | .862 | .823 | .769 | .837 | .861 | .877 |
| 40 | .844 | .796 | | .799 | .839 | .853 |
| 50 | .814 | .754 | | .717 | .810 | .845 |

(a) Accuracy of qualitative derivatives.

| | LWR | Pairs | $\tau$-regression $\kappa = 30$ | 50 | 70 | 100 |
|---|---|---|---|---|---|---|
| $k = 10$ | .509 | .519 | .516 | .512 | .510 | .509 |
| 20 | .515 | .531 | .509 | .518 | .522 | .510 |
| 30 | .515 | .523 | .510 | .513 | .514 | .515 |
| 40 | .521 | .555 | | .511 | .507 | .511 |
| 50 | .516 | .583 | | .507 | .508 | .515 |

(b) Accuracy of qualitative models induced by C4.5.

Table 3: Results of experiments with $f(x, y) = \sin x \sin y$.

Apart from the failure of $\tau$-regression on $\partial F / \partial m$, which we cannot explain, the proportion of correctly computed qualitative partial derivatives (Table 4a) is around 100%. C4.5 correctly induced single-node trees predicting $F = Q(+m)$, $F = Q(+v)$ and $F = Q(-r)$.

*Gravitation.* Newton's law of universal gravitation states that a point mass attracts another point mass by a force pointing along the line intersecting both points with the force equal to

$$F = G\frac{m_1 m_2}{r^2},$$

where $F$ is gravitational force, $G$ is the gravitational constant, $m_1$ and $m_2$ are the points' masses, and $r$ is the distance between the two point masses. We prepared a data set where $m_1$ and $m_2$ were sampled randomly from $[0.1, 1]$ and $r$ from $[0.1, 2]$.

Padé again reached almost 100% accuracy (Table 4b) and C4.5 induced correct models $F = Q(+m_1)$, $F = Q(+m_2)$ and $F = Q(-r)$.

*Pendulum.* The period $T$ of the first swing of a pendulum w.r.t. its length

| variable | $m$ | $v$ | $r$ |
|---|---|---|---|
| LWR | 1.00 | 1.00 | 1.00 |
| $\tau$-regression | .86 | .97 | .94 |
| parallel pairs | 1.00 | 1.00 | .98 |

(a) Centripetal force

| variable | $m_1$ | $m_2$ | $r$ |
|---|---|---|---|
| LWR | .964 | .96 | 1.00 |
| $\tau$-regression | .88 | .87 | .99 |
| parallel pairs | .96 | .96 | 1.00 |

(b) Gravitational force

| variable | $l$ | $\phi$ |
|---|---|---|
| LWR | 1.00 | .98 |
| $\tau$-regression | 1.00 | .83 |
| parallel pairs | 1.00 | .97 |

(c) Pendulum

Table 4: Accuracy of computed partial derivatives for domains from physics.

$l$, the mass of the bob $m$ and the initial angle $\varphi$ is[2]

$$T = 2\pi\sqrt{\frac{l}{g}}(1 + \frac{1}{16}\varphi_0^2 + \frac{11}{3072}\varphi_0^4 + \frac{173}{737280}\varphi_0^6 + \ldots).$$

The mass was again chosen randomly from $[0.1, 1]$, the length was from $[0.1, 2]$, the angle was between $-180$ and $180$ degrees, and $g = 9.81$.

Table 4c shows the accuracies for dependency of the force on the length and the angle. All methods achieved $100\%$ accuracy for the length, *i.e.* for each learning example the qualitative derivative equals the ground truth.

The qualitative model describing the relation between $T$ and $\varphi$ depends on the value of $\varphi$. C4.5 induced a qualitative tree shown in Fig. 2 from data labelled by LWR. The models based on $\tau$-regression and parallel pairs are the same – the only difference is in the split threshold at the root node which is $-2.7\,\mathrm{deg}$ for $\tau$-regression and $-2.28\,\mathrm{deg}$ for parallel pairs. The tree states that for negative

---

[2]http://en.wikipedia.org/wiki/Pendulum

Figure 2: The qualitative model describing the relation between $T$ and $\varphi$.



(a) Data set $x^2 - y^2$.
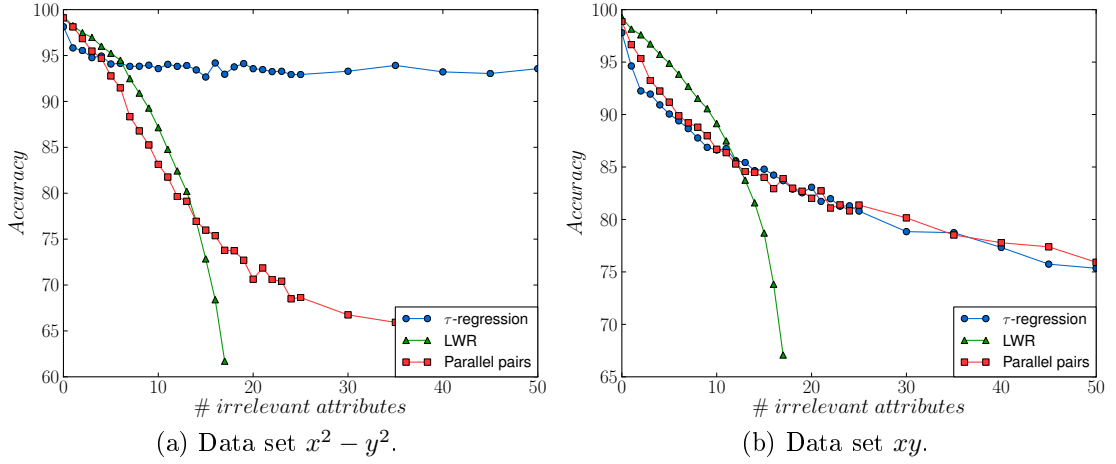
(b) Data set $xy$.

Figure 3: Accuracy for different number of additional irrelevant attributes.

angles the period $T$ decreases with $\varphi$ while for positive angles it increases with $\varphi$. In other words, the period increases with the absolute value of $\varphi$.

*4.2. Effect of irrelevant attributes*

Many real-world data sets include large number of attributes with no effect on the function value. The following experiment shows the robustness of Padé in such cases.

We consider data sets obtained by sampling $f(x, y) = x^2 - y^2$ and $f(x, y) = xy$ as described above and observe the correctness of partial derivatives computed by Padé when 0–50 random attributes from the same interval as $x$ and $y$, $[-10, 10]$, are added to the domain. The graphs in Fig. 3 show how the accuracy drops with the increasing number of added irrelevant attributes for each method.

|              | $\sigma = 0$ | $\sigma = 10$ | $\sigma = 30$ | $\sigma = 50$ |
|--------------|------|------|------|------|
| LWR          | .993 | .962 | .878 | .795 |
| $\tau$-regression | .981 | .945 | .848 | .760 |
| parallel pairs | .992 | .924 | .771 | .680 |

(a) Correctness of computed derivatives

|              | $\sigma = 0$ | $\sigma = 10$ | $\sigma = 30$ | $\sigma = 50$ |
|--------------|------|------|------|------|
| LWR          | .996 | .978 | .956 | .922 |
| $\tau$-regression | .991 | .976 | .949 | .917 |
| parallel pairs | .995 | .966 | .949 | .901 |

(b) Correctness of qualitative models induced by C4.5

Table 5: Effect of noise on the accuracy of computed qualitative partial derivatives for $f(x, y) = x^2 - y^2$ with random noise from $N(0, \sigma)$.

We observe a steep drop in the accuracy of LWR which we can ascribe to the fact that LWR is multivariate and is solving an ill-conditioned system as the number of attributes approaches $k$. Parallel pairs and $\tau$-regression avoid this problem by computing the derivative by a single attribute at a time.

As for comparison between $\tau$-regression and parallel pairs, we performed the same experiment with other similar functions and found that no method consistently outperforms the other.

*4.3. Coping with noise*

In this experiment we add various amounts of noise to the function value. The target function is $f(x, y) = x^2 - y^2$ defined on $[-10, 10] \times [-10, 10]$ which puts $f$ in $[-100, 100]$. We added Gaussian noise with a mean of 0 and variance 0, 10, 30, and 50, *i.e.* from no noise to the noise in the range comparable to the signal itself. We measured the accuracy of derivatives and models induced by C4.5. Since the data contains noise, we set the C4.5's parameter $m$ (minimal number of examples in a leaf) to 10% of the examples of our data set, $m = 100$. We repeated the experiment with each amount of noise 100 times and computed the average accuracies.

The results are shown in Table 5. Padé is quite robust despite the huge amount of noise. As in other experiments on artificial data sets, we again observed that C4.5 is able to learn almost perfect models despite the drop in correctness of derivatives at higher noise levels.

*4.4. Stability on real-world data sets*

In measurements of stability we ran Padé on six UCI ML repository [32] data sets (servo, imports-85, galaxy, prostate, housing, auto-mpg) and the four artificial data sets $(xy, x^2 - y^2, x^3 - y, \sin x \sin y)$ which were already used in previous experiments. For each data set we created 100 bootstrap samples and for each sample we computed qualitative partial derivatives of the class w.r.t. all continuous attributes for each example. We defined stability w.r.t. attribute $x_i$ for example $\mathbf{x_m}$, $\beta_{\mathbf{x_m}}(x_i)$, as the proportion of the prevailing derivatives $(Q_{\mathbf{x_m}}(+x_i)$ or $Q_{\mathbf{x_m}}(-x_i))$ among all positive and negative derivatives predicted for $\mathbf{x_m}$,

$$\beta_{\mathbf{x_m}}(x_i) = \frac{max\left(\#Q_{\mathbf{x_m}}(+x_i), \#Q_{\mathbf{x_m}}(-x_i)\right)}{\#Q_{\mathbf{x_m}}(+x_i) + \#Q_{\mathbf{x_m}}(-x_i)}. \tag{17}$$

The total stability with regard to attribute $x_i$ is the average stability across all examples in the data set. We prefer $\beta$ close to 1, which happens when most of the qualitative predictions for this derivative are the same. The lowest possible $\beta$ is 0.5.

The settings for all methods were $k = 20$ and $\kappa = 50$ for $\tau$-regression. The lower value of $\kappa$ was used due to a smaller number of examples. Results in Table 6 show the average stability of the partial derivative of the outcome w.r.t. to each attribute.

Altogether there are 46 continuous attributes in all domains. We ranked performances of methods for each attribute, where rank 1 was assigned to the best method and rank 3 to the worst one. Averaging over all attributes showed that the parallel pairs method is the most stable with a rank of 1.80, slightly behind is the $\tau$-regression method with 1.84, while the LWR method achieved a considerably worse rank of 2.35.

## 5. Case study: Billiards

To show a practical use of the proposed methods, we analyse a problem from billiards. Billiards is a common name for table games played with a stick and a set of balls, such as snooker or pool and their variants. The goals of the games vary, but the main idea is common to all: the player uses the stick to stroke the cue ball aiming at another ball to achieve the desired effect. The friction between the table and the balls, the spin of the cue ball and the collision of the balls combine into a very complex physical system [33]. However, an amateur player can learn the basic principles of how to stroke the cue ball without knowing much about the physics behind it.

| data set/attribute | LWR | $\tau$ | PP |
|---|---|---|---|
| auto | | | |
| displacement | .76 | .80 | .75 |
| horsepower | .71 | .84 | .75 |
| weight | .77 | .81 | .78 |
| acceleration | .75 | .76 | .77 |
| imports-85 | | | |
| normalized-losses | .68 | .76 | .72 |
| wheel-base | .66 | .85 | .89 |
| length | .72 | .85 | .82 |
| width | .67 | .82 | .94 |
| curb-weight | .88 | .88 | .96 |
| height | .72 | .79 | .75 |
| engine-size | .71 | .80 | .91 |
| bore | .67 | .78 | .85 |
| compression-ratio | .68 | .81 | .79 |
| stroke | .65 | .80 | .81 |
| horsepower | .64 | .89 | .93 |
| peak-rpm | .66 | .83 | .75 |
| highway-mpg | .66 | .87 | .84 |
| city-mpg | .65 | .87 | .88 |
| housing | | | |
| crim | .71 | .80 | .78 |
| indus | .68 | .85 | .81 |
| nox | .70 | .80 | .85 |
| rm | .87 | .88 | .92 |
| age | .82 | .79 | .79 |
| dis | .71 | .81 | .81 |
| tax | .70 | .88 | .84 |
| ptratio | .67 | .85 | .84 |
| b | .68 | .76 | .75 |
| lstat | .71 | .91 | .93 |

| data set/attribute | LWR | $\tau$ | PP |
|---|---|---|---|
| galaxy | | | |
| east.west | .93 | .90 | .88 |
| north.south | .93 | .90 | .80 |
| radial.position | .93 | .90 | .83 |
| prostate | | | |
| lcavol | .92 | .82 | .93 |
| lweight | .78 | .80 | .80 |
| age | .69 | .79 | .79 |
| lbph | .69 | .78 | .78 |
| lcp | .71 | .77 | .71 |
| servo | | | |
| pgain | .95 | .94 | 1.00 |
| vgain | .88 | .72 | .85 |
| $xy$ | | | |
| x | .99 | .95 | .99 |
| y | .99 | .94 | .98 |
| $x^3 - y$ | | | |
| x | 1.00 | 1.00 | 1.00 |
| y | .77 | .75 | .78 |
| $x^2 - y^2$ | | | |
| x | .99 | .95 | .99 |
| y | .99 | .96 | .99 |
| $\sin x \sin y$ | | | |
| x | .91 | .88 | .86 |
| y | .91 | .89 | .86 |

Table 6: Stability on real and artificial data sets. The numbers correspond to stabilities of LWR, $\tau$-regression and parallel pairs, respectively.

(a) A stroke with azimuth 0. The cue ball hits the black one with a *full hit*.

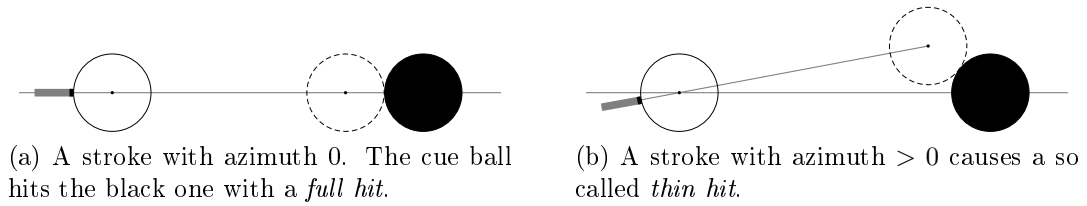(b) A stroke with azimuth > 0 causes a so called *thin hit*.

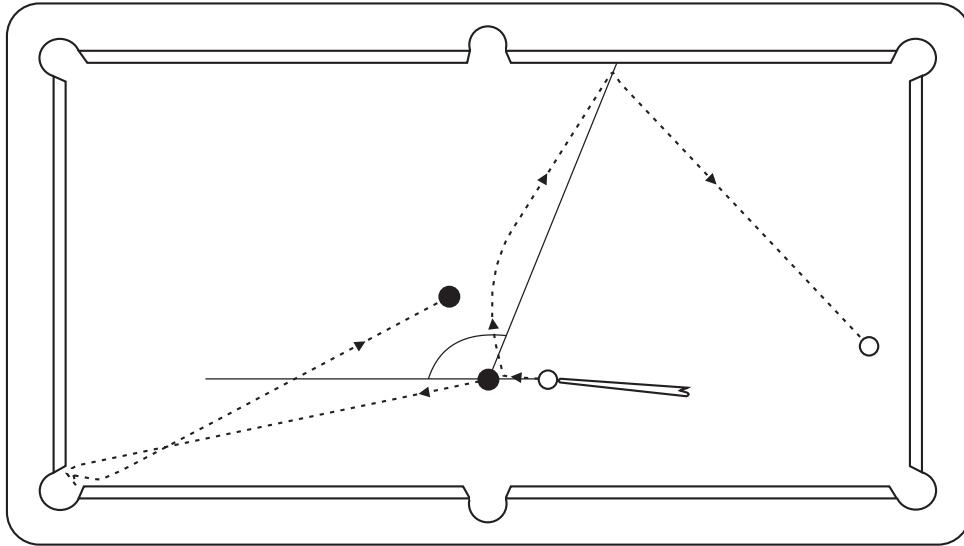Figure 4: Strokes with different fullness of hit (azimuth).



Figure 5: Example trace from the simulator.

Our goal is to induce a qualitative model describing the relation between the azimuth (degree of fullness of hit; see Fig. 4) and the reflection angle of the cue ball after the collision with another ball. We define the reflection angle as the angle between the stroke direction and the line connecting the positions of the cue ball's collision with the black ball and its next collision, either with the cushion or with the black ball again (Fig. 5).

We model the relation between the azimuth and the angle as it depends on four attributes (Table 7): the horizontal and the vertical angle of the stick, called *azimuth* and *elevation*, respectively, as well as the *follow* (forward/backward spin; see Fig. 6) of the cue ball and the *velocity* of the stroke.

Figure 6: Spinning the cue ball with a follow (above centre) or a draw (below centre). Dots represent the point at which the cue ball is hit with the stick.

| name | range | description |
|---|---|---|
| azimuth | [-15, 15] *deg* | the horizontal angle of the stick |
| elevation | [0, 30] *deg* | the vertical angle of the stick |
| velocity | [3, 5] *m/s* | velocity of the stick |
| follow | [−.5r, .5r] | forward/backward spin caused by hitting the ball above/below the centre |
| angle | [-180, 180) *deg* | the reflection angle of the cue ball |

Table 7: The observed variables in the billiards case study.

We used the billiards simulator [34] to collect a sample of 5000 randomly chosen scenarios. We chose to use $\tau$-regression with parameters as usual, $\kappa = 100$, $k = 20$. We induced a qualitative tree using C4.5 and asked the experts for their interpretation of the model.
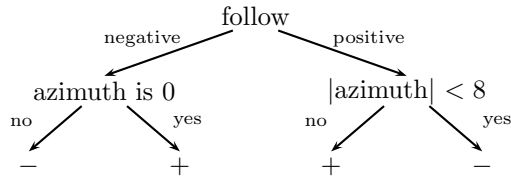
The tree induced by C4.5 is shown in Fig. 7a. The experts found the model much easier to understand than a set of equations and the knowledge easier to transfer to the actual game. They were also able to recast the original tree into a simpler model (Fig. 7b). As such, it can serve conceptualization of the domain useful for beginners.

The tree first splits on attribute *follow*. For negative values (the cue ball is hit below the centre), the backward spin results in inversely proportional relation between the reflection angle and the *azimuth*, *i.e.*, increasing the *azimuth* decreases the *angle* and vice versa. The exception at the *azimuth* of zero is the artefact of the coordinate system: the reflection angle at this *azimuth* crosses the line of discontinuity at $\pm 180$ deg, so an increase from, say, 179 to 181 degrees is recorded as a decrease from 179 to -179 degrees.

Positive values of *follow* give the cue ball an extra amount of forward spin,

(a) The original tree induced by C4.5



(b) Reinterpretation by domain experts

Figure 7: The qualitative model describing the relation between the reflection angle and azimuth.

which adds to the spin that the cue ball acquires due to rolling. In this case, *azimuth* has a different effect on the black ball. Having a forward spin, the cue ball has a tendency to roll forward after the collision. For *azimuth* approximately between $-8$ and $8$ degrees, increasing (decreasing) the *azimuth* also increases (decreases) the reflection angle. For greater absolute values of *azimuth* the cue ball only partially hits the black ball which causes the reflection angle to become negatively related with the *azimuth*.

Our model also correctly recognized the two important attributes, *follow* and *azimuth*. The other two, *elevation* and *velocity* indeed have no effect on reflection angle in our context, according to the expert opinion.

We have consequently tested the obtained qualitative relations in the simulator by simulating small changes in azimuth while leaving the values of other attributes constant and observing the reflection angle. We found all rules correct apart from some small deviations regarding the numerical values of the splits in the qualitative tree.

## 6. Discussion and conclusion

The three methods for computation of derivatives differ in their theoretical background, which leads to differences in their performance.

LWR is a multivariate method, computing all derivatives at once, while the other two are univariate. This mostly puts LWR at a disadvantage. LWR is more sensitive to the size of data set: if the number of attributes is comparable to (or even greater) than the number of points on which the regression is computed (in our case, this is the number of examples in the neighbourhood, $k$, and not the total size of the data set), the LWR is solving an ill-conditioned system. This problem clearly surfaces in our experiments with irrelevant attributes. The other two methods take one attribute at a time and do not run into this problem. When the number of attributes is small, LWR however gives satisfactory results.

LWR also fared considerably worse on real-world data sets where we evaluated the methods indirectly by testing their stability. The lower stability of LWR reflects its higher complexity. Based on multivariate regression, the number of parameters that LWR optimises equals the number of attributes. The higher complexity of the LWR method allows it to fit the data better, *i.e.* it has lower bias at the cost of higher variance. Therefore, LWR can in theory better estimate derivatives when compared to the other two, but it can overfit on data with many attributes. This property is also manifested in our experiments, where LWR performed considerably worse on domains with many attributes (housing and imports-85), while it scored best on the five domains which only have two attributes.

Being univariate, the $\tau$-regression and the parallel pairs methods have to select examples lying in the direction of differentiation to exclude the influence of other attributes on function value. To gather the same number of examples as the LWR, they expand the neighbourhood. The $\tau$-regression does so by extending the neighbourhood in direction of the axis of differentiation and the parallel pairs method differentiates along lines parallel to the actual axis of differentiation. The accuracy of computed derivatives should profit from excluding other attributes and suffer from extending the neighbourhood. The effect of this is visible on target functions $x^3 - y$ and $\sin x \sin y$. The former tests how well the method copes with effects of $x$ when differentiating by $y$; LWR fails here, while parallel pairs work well. The latter concept requires a small neighbourhood, so LWR fares better here, although the difference in performance is not as large as in for the former concept.

The difference between the $\tau$-regression and parallel pairs is subtle. The $\tau$-

regression computes the derivative inside a somewhat longer tube, whose length is governed by the parameter $\kappa$, which we typically set to 50–100. It assumes that the function is linear in a longer region, and the error comes mostly from violating this assumption. Parallel pairs use a smaller vicinity; in most experiments we set $\kappa$ and $k$ to 20. The function is thus required to be linear only close to $\mathbf{x_0}$. The price to pay is that the derivative is smoothed over a hyper-disc perpendicular to the axis of differentiation instead of being computed at $\mathbf{x_0}$, as in $\tau$-regression. This may cause problems if the derivative strongly depends upon other arguments of the function. We were however unable to experimentally find any practical consequences of this theoretical difference between the $\tau$-regression and parallel pairs. In general, recommending one method for a particular context is equally impossible as recommending a machine learning algorithm for a certain domain. A general approach should be to try both.

Experiments demonstrate that while computation of partial derivatives is often difficult and prone to errors, even quite unreliable estimates usually result in correct qualitative models. Even a modest learning algorithm such as C4.5 is able to reduce the noise in the computed derivatives and build models which almost perfectly match the target function.

Apart from ignoring the noisy labels, the C4.5 task was rather trivial for most artificial domains as it required inducing a tree with one or two leaves. It was put to a more serious test in a case study from billiards. It has induced a correct model which also demonstrates the basic advantage of qualitative modelling over standard regression modelling: the model we induced presents the general principles which can be used by a player, while a regression model would describe complex numerical relationships which would probably be difficult to comprehend and certainly impossible to use at the billiards table.

While this model is still rather small and simple, real-world qualitative models induced by experts, such as Samuelson [3] and Jeffries and May [1, 2], do not get any more complex than this neither. It may be that either the world is simpler than expected when described qualitatively, or that humans are used to reason using simple models which skip over the details, which are included only in specific models for particular contexts.

In summary, we introduced a new approach to learning qualitative models which differs from existing approaches by its trick of translating the learning problem into a classification problem and then applying the general-purpose learning methods to solve it. To focus the paper, we mostly explored the first step which involves the estimation of partial derivatives. The second step opens a number of other interesting research problems not explicitly discussed here.

One is exploration of other, more powerful learning algorithms. Is our above assumption of the simplicity of the world in qualitative terms correct, or would using, say, support vector machines, result in much better models? Another interesting question is how to combine models for derivatives with respect to individual arguments into a single model describing the function's behaviour with respect to all arguments. One approach may be to use multi-label learning methods to model all derivatives at once. We leave these questions open for further research in the area.

**Acknowledgement**

[1] C. Jeffries, Qualitative stability and digraphs in model ecosystems, Ecology 55 (6) (1974) 1415–1419.

[2] R. M. May, Qualitative stability in model ecosystems, Ecology 54 (3) (1973) 638–641.

[3] P. A. Samuelson, Foundations of Economic Analysis, Harvard University Press; Enlarged edition, 1983.

[4] K. Forbus, Qualitative process theory, Artificial Intelligence 24 (1984) 85–168.

[5] J. de Kleer, J. S. Brown, A qualitative physics based on confluences, Artificial Intelligence 24 (1984) 7–83.

[6] B. Kuipers, Qualitative simulation, Artificial Intelligence 29 (1986) 289–338.

[7] J. Kalagnanam, H. A. Simon, Y. Iwasaki, The mathematical bases for qualitative reasoning, IEEE Intelligent Systems 6 (2) (1991) 11–19.

[8] J. R. Kalagnanam, Qualitative analysis of system behaviour, Ph.D. thesis, Pittsburgh, PA, USA (1992).

[9] J. Kalagnanam, H. A. Simon, Directions for qualitative reasoning, Computational Intelligence 8 (2) (1992) 308–315.

[10] E. Coiera, Generating qualitative models from example behaviours, Tech. Rep. 8901, University of New South Wales (1989).

[11] D. Hau, E. Coiera, Learning qualitative models of dynamic systems, Machine Learning Journal 26 (1997) 177–211.

[12] S. Ramachandran, R. Mooney, B. Kuipers, Learning qualitative models for systems with multiple operating regions, in: Working Papers of the 8th International Workshop on Qualitative Reasoning about Physical Systems, Japan, 1994.

[13] B. Richards, I. Kraan, B. Kuipers, Automatic abduction of qualitative models, in: Proceedings of the National Conference on Artificial Intelligence, AAAI/MIT Press, 1992.

[14] A. Say, S. Kuru, Qualitative system identification: deriving structure from behavior, Artificial Intelligence 83 (1996) 75–141.

[15] I. Bratko, S. Muggleton, A. Varšek, Learning qualitative models of dynamic systems, in: Proc. Inductive Logic Programming ILP-91, 1991, pp. 207–224.

[16] S. M. Coghill, G. M.; Garrett, R. D. King, Learning qualitative models in the presence of noise, in: Proc. Qualitative Reasoning Workshop QR'02, 2002.

[17] G. M. Coghill, A. Srinivasan, R. D. King, Qualitative system identification from imperfect data, J. Artif. Int. Res. 32 (1) (2008) 825–877.

[18] S. Džeroski, L. Todorovski, Discovering dynamics: from inductive logic programming to machine discovery, Journal of Intelligent Information Systems 4 (1995) 89–108.

[19] S. Džeroski, L. Todorovski, Discovering dynamics, in: International Conference on Machine Learning, 1993, pp. 97–103.

[20] L. Todorovski, Using domain knowledge for automated modeling of dynamic systems with equation discovery, Ph.D. thesis, University of Ljubljana, Ljubljana, Slovenia (2003).

[21] H. Kay, B. Rinner, B. Kuipers, Semi-quantitative system identification, Artificial Intelligence 119 (2000) 103–140.

[22] R. K. Gerçeker, A. Say, Using polynomial approximations to discover qualitative models, in: Proc. of the 20th International Workshop on Qualitative Reasoning, Hanover, New Hampshire, 2006.

[23] I. Bratko, I. Mozetič, N. Lavrač, KARDIO: a study in deep and qualitative knowledge for expert systems, MIT Press, Cambridge, MA, USA, 1989.

[24] I. Mozetič, Learning of qualitative models, in: EWSL, 1987, pp. 201–217.

[25] I. Mozetič, The role of abstractions in learning qualitative models, in: Proc. Fourth Int. Workshop on Machine Learning, Morgan Kaufmann, 1987.

[26] D. Šuc, I. Bratko, Induction of qualitative trees, in: L. De Raedt, P. Flach (Eds.), Proceedings of the 12th European Conference on Machine Learning, Springer, 2001, pp. 442–453, Freiburg, Germany.

[27] D. Šuc, Machine Reconstruction of Human Control Strategies, Vol. 99 of Frontiers in Artificial Intelligence and Applications, IOS Press, Amsterdam, The Netherlands, 2003.

[28] I. Bratko, D. Šuc, Learning qualitative models, AI Magazine 24 (4) (2003) 107–119.

[29] C. Atkeson, A. Moore, S. Schaal, Locally weighted learning, Artificial Intelligence Review 11 (1997) 11–73.

[30] S. Ben-David, U. von Luxburg, D. Pal, A sober look at clustering stability, in: 19th Annual Conference on Learning Theory, Springer, Berlin, Germany, 2006, pp. 5–19.

[31] J. Demšar, B. Zupan, G. Leban, T. Curk, Orange: From experimental machine learning to interactive data mining, in: Proceedings of PKDD 2004, 2004, pp. 537–539.

[32] A. Asuncion, D. J. Newman, UCI machine learning repository (2007).

[33] D. G. Alciatore, The Illustrated Principles of Pool and Billiards, Sterling; 1st edition, 2004.

[34] D. Papavasiliou, Billiards manual, Tech. rep. (2009).
     URL http://www.nongnu.org/billiards/