# Metadata of the chapter that will be visualized online

| | | |
|---|---|---|
| Chapter Title | Trading Volatility Using Highly Accurate Symbolic Regression | |
| Copyright Year | 2015 | |
| Copyright Holder | Springer International Publishing Switzerland | |
| Corresponding Author | Family Name | **Korns** |
| | Particle | |
| | Given Name | **Michael F.** |
| | Suffix | |
| | Organization | |
| | Address | 2240 Village Walk Drive, Henderson, NV, 89052, USA |
| | Email | mkorns@korns.com |

| Abstract | Research efforts, directed at increasing the accuracy and dependability of Symbolic Regression (SR), have resulted in significant improvements in symbolic regression's range, accuracy, and dependability. Previous research has also demonstrated the practicability of estimating corporate forward 12 month earnings, using advanced symbolic regression. In this paper we put these prior results and techniques together to select a 100 stock semi-passive index portfolio, extracted from the Value Line Timeliness stocks (*Value Line*), which delivers consistent performance in both bull and bear decades and we will compare its performance to the Standard & Poors 100 index.

We intend to produce our 100 stock semi-passive index buy list on a weekly basis using automated forward 12 month EPS (*ftmEPS*) prediction involving the analysis of many securities, involving multiple training regressions each on hundreds of thousands of training examples. Plus the timeliness issue will require that our analytic tools be strong and thoroughly matured. The 100 stock buy list will be the foundation for a new semi-passive Value Line 100 index fund which should have great appeal to many high net worth clients, enjoy low management costs, and be easily acceptable to the compliance and regulatory authorities.

Valuation of Value Line securities via their forward 12 month price earnings ratio (*ftmPE*) is a very common securities valuation method in the industry. Obviously the *ftmPE* valuation depends heavily on the estimate of forward 12 month corporate earnings per share (*ftmEPS*). Several obvious inputs to the *ftmEPS* prediction process are the past earnings time series plus one or more analyst predictions.

Valuation via *ftmEPS* is a necessary but not a sufficient attraction for a semi-passive index fund. So we will introduce the advantages of trading volatility. Our thesis will be that emotional trading patterns tend to make markets less efficient.

The efficient market hypothesis depends upon equal access to information and rational trading patterns. Trading on insider information is illegal in most developed securities markets; however, trading when others are emotional is unregulated. In this paper we will develop a set |

of factors—all of which incorporate a measure of volatility indicating possible overly emotional trading patterns. The theme of our new semi-passive index fund will be "*Buy value from those who are selling in a highly emotional state*".

# Chapter 21
# Trading Volatility Using Highly Accurate Symbolic Regression

**Michael F. Korns**

## 21.1 Introduction

The discipline of Symbolic Regression (SR) has matured significantly in the last few years. There is at least one commercial package on the market for several years (http://www.rmltech.com/). There is now at least one well documented commercial symbolic regression package available for Mathematica (www.evolved-analytics. com). There is at least one very well done open source symbolic regression package available for free download (http://ccsl.mae.cornell.edu/eureqa).

In addition to our own ARC system (Korns 2013, 2014), currently used internally for massive financial data nonlinear regressions, there are a number of other mature symbolic regression packages currently used in industry including Smits et al. (2010) and Castillo et al. (2010). Plus there is an interesting work in progress by McConaghy et al. (2009).

Research efforts, directed at increasing the accuracy and dependability of Symbolic Regression (SR), have resulted in significant improvements in symbolic regression's range, accuracy, and dependability (Korns 2013, 2014). Previous research has also demonstrated the practicability of estimating corporate forward 12 month earnings, using advanced symbolic regression (Korns 2012a, b). In this paper we put these prior results and techniques together to select a 100 stock semi-passive index portfolio (*VEP100*), from the Value Line Timeliness (*Value Line*), which delivers consistent performance in both bull and bear decades.

We intend to produce our VEP100 buy list on a weekly basis using automated *ftmEPS* prediction involving the analysis of many securities, involving multiple training regressions each on hundreds of thousands of training examples. Plus the timeliness issue will require that our analytic tools be strong and thoroughly

M.F. Korns (✉)
2240 Village Walk Drive, Henderson, NV 89052, USA
e-mail: mkorns@korns.com

matured. Our new VEP100 semi-passive index fund should have great appeal to many high net worth clients, enjoy low management costs, and be easily acceptable to the compliance and regulatory authorities.

Valuation of securities via their forward 12 month price earnings ratio (*ftmPE*) is a very common securities valuation method in the industry. Obviously the *ftmPE* valuation depends heavily on the estimate of forward 12 month corporate earnings per share (*ftmEPS*). Obvious inputs to the *ftmEPS* prediction process are the past earnings time series plus one or more analyst predictions.

Valuation via *ftmEPS* is a necessary but not a sufficient attraction for a semi-passive index fund. So we will introduce the advantages of trading volatility. Our thesis will be that emotional trading patterns tend to make markets less efficient.

The efficient market hypothesis assumes rational trading patterns and equal and open access to information. Trading on insider information is illegal in most developed securities markets; but, *trading when others are emotional is unregulated*. In this paper we will develop a set of factors—all of which incorporate a measure of volatility indicating possible overly emotional trading patterns. The theme of our new VEP100 semi-passive index fund will be "*Buy value from those who are selling in a highly emotional state*".

Now would be a good time to provide an overview general introduction to symbolic regression as follows.

Symbolic Regression is an approach to general nonlinear regression which is the subject of many scholarly articles in the Genetic Programming community. A broad generalization of general nonlinear regression is embodied as the class of *Generalized Linear Models* (GLMs) as described in Nelder and Wedderburn (1972). A GLM is a linear combination of **I** basis functions $\mathbf{B_i}$; $i = 1,2, \ldots, I$, a dependent variable **y**, and an independent data point with **M** features $\mathbf{x} = <x_1, x_2, x_3, \ldots, x_m>$: such that

$$y = \gamma(x) = C_0 + \sum_{i=1}^{I} c_i B_i(x) + err \qquad (21.1)$$

As a broad generalization, GLMs can represent any possible nonlinear formula. However the format of the GLM makes it amenable to existing linear regression theory and tools since the GLM model is linear on each of the basis functions $\mathbf{B_i}$.

For a given vector of dependent variables, Y, and a vector of independent data points, X, symbolic regression will search for a set of basis functions and coefficients which minimize *err*. In Koza (1992) the basis functions selected by symbolic regression will be formulas as in the following examples:

$$B_1 = x_3 \qquad (21.2)$$

$$B_2 = x_1 + x_4 \qquad (21.3)$$

$$B_3 = \text{sqrt}(x_2) / \tan(x_5/4.56) \qquad (21.4)$$

$$B_4 = \tanh\left(\cos\left(x_2 * .2\right) * \text{cube}\left(x_5 + \text{abs}\left(x_1\right)\right)\right) \tag{21.5}$$

If we are minimizing the least squared error, **LSE**, once a suitable set of basis    63
functions {**B**} have been selected, we can discover the proper set of coefficients    64
{**C**} deterministically using standard univariate or multivariate regression. The value    65
of the GLM model is that one can use standard regression techniques and theory.    66
Viewing the problem in this fashion, we gain an important insight. Symbolic    67
regression does not add anything to the standard techniques of regression. The value    68
added by symbolic regression lies in its abilities as a search technique: how quickly    69
and how accurately can SR find an optimal set of basis functions {**B**}.    70

The immense size of the search space provides ample need for improved search    71
techniques In standard Koza-style tree-based Genetic Programming (Koza 1992)    72
the genome and the individual are the same Lisp s-expression which is usually    73
illustrated as a tree. Of course the tree-view of an s-expression is a visual aid,    74
since a Lisp s-expression is normally a list which is a special Lisp data structure.    75
Without altering or restricting standard tree-based GP in any way, we can view the    76
individuals not as trees but instead as s-expressions such as this depth 2 binary tree    77
s-exp: (/ (+x_2 3.45) (*x_0 x_2)), or this depth 2 irregular tree s-exp: (/ (+x_2 3.45) 2.0).    78

In standard GP, applied to symbolic regression, the non-terminal nodes are all    79
operators (implemented as Lisp function calls), and the terminal nodes are always    80
either real number constants or features. The maximum depth of a GP individual is    81
limited by the available computational resources; but, it is standard practice to limit    82
the maximum depth of a GP individual to some manageable limit at the start of a    83
symbolic regression run.    84

Given any selected maximum depth k, it is an easy process to construct a maximal    85
binary tree s-expression $U_k$, which can be produced by the GP system without    86
violating the selected maximum depth limit. As long as we are reminded that each f    87
represents a function node while each t represents a terminal node, the construction    88
algorithm is simple and recursive as follows.    89

$$U_0 : t$$

$$U_1 : (f\ t\ t)$$

$$U_2 : (f\ (f\ t\ t)\ (f\ t\ t))$$

$$U_3 : (f\ (f\ (f\ t\ t)\ (f\ t\ t))\ (f\ (f\ t\ t)\ (f\ t\ t)))$$

$$U_k : (f\ U_{k\text{-}1} U_{k\text{-}1})$$

Any basis function produced by the standard GP system will be represented by at 90
least one element of $U_k$. In fact, $U_k$ is isomorphic to the set of all possible basis 91
functions generated by the standard GP system. 92

Given this formalism of the search space, it is easy to compute the size of the 93
search space, and it is easy to see that the search space is huge even for rather simple 94
basis functions. For our use in this chapter the function set will be the following 95
functions: $\mathbf{F} = \{+ - * / $ **abs sqrt square cube cos sin tan tan h log exp max** 96
**min$\aleph$**$\}$ (where $\aleph(a,b) = \aleph(a) = a$). The terminal set is the features $\mathbf{x_0}$ thru $\mathbf{x_m}$ and 97
the real constant $\mathbf{c}$, which we shall consider to be $2^{64}$ in size. Where $|\mathbf{F}| = 17$, 98
$\mathbf{M}=20$, and $\mathbf{k} = 0$ , the search space is $S_0 = \mathbf{M} + 2^{64} = 20 + 2^{64} = 1.84 \times 10^{19}$. 99
Where $\mathbf{k} = 1$, the search space is $S_1 = |\mathbf{F}| * S_0 * S_0 = 5.78 \times 10^{39}$. Where $\mathbf{k} = 2$, 100
the search space grows to $S_2 = |\mathbf{F}| * S_1 * S_1 = 5.68 \times 10^{80}$. For $\mathbf{k} = 3$, the search 101
space grows to $S_3 = |\mathbf{F}| * S_2 * S_2 = 5.5 \times 10^{162}$. Finally if we allow three basis 102
functions $\mathbf{B} = 3$ for financial applications, then the final size of the search space 103
is $S_3 * S_3 * S_3 = 5.5 \times 10^{486}$. 104

## 21.2 Methodology 105

Creating the weekly buy list for a modern semi-passive index fund requires many 106
fully automated multiple regressions, all of which must be run in a timely fashion, 107
and all of which must fit together seamlessly without human intervention. Our 108
methodology is influenced by the practical issues of applying symbolic regression 109
to the real world investment finance environment. First there is the issue that form 110
of each symbolic regression must be preapproved by the regulatory authorities, the 111
compliance officer, management, and clients. Second there is the issue of adapting 112
symbolic regression to a real world financial application with massive 113
amounts of data. Third there is the issue of modifying symbolic regression, as 114
practiced in academia, to conform to the very difficult U.S. Securities Exchange 115
Commission regulatory compliance environment. 116

Weekly preparation of our VEP100 semi-passive index fund buy list will require 117
∼1502 fully automated regressions (*as many as there are Value Line Timeliness* 118
*stocks that week*). For each of the ∼1500 Value Line Timeliness stocks, a set of 119
pre-approved earnings estimate inputs will be fed into a multiple linear regression 120
for each stock, resulting in an interim forward 12 month earnings per share estimate 121
for the stock. This will require ∼1500 regressions; but, they are relatively quick 122
multiple linear regressions. Next, a set of preapproved earnings estimate inputs 123
plus the interim *ftmEPS* estimate produced by the linear regressions will be input 124
to a nonlinear weighted regression on all ∼1500 stocks. This expensive nonlinear 125
weighted regression will produce a final *ftmEPS* estimate for each of the Value 126
Line stocks. Finally, a set of preapproved z Score factor inputs plus the interim 127
*ftmEPS* estimate produced by the linear regressions will be input to a nonlinear 128
logistic regression on all ∼1500 stocks. This final very expensive nonlinear logistic 129
regression will produce a final *expected forward 12 month total return* estimate for 130
each of the Value Line stocks. 131

We use only statistical best practices out-of-sample testing methodology. For each regression, a matrix of independent variables will be constructed solely from the prior 10 years of historical data—520 weeks. No forward looking data will be allowed. This is very important because it will be the subject of detailed regulatory due diligence reviews. Then the preapproved regression model will be applied to produce the dependent variable.

For the forward estimation of corporate earnings, this paper uses an historical database of the Value Line stocks with daily price and volume data, weekly analyst estimates, and quarterly financial data from January 1990 to the December 2009. The data has been assembled from reports published at the time, so the database is highly representative of what information was realistically available at the point when trading decisions were actually made. No forward looking data is included in any historical point in the database.

From all of this historical data, 20 years (*1990 thru 2009*) have been used to produce the results shown in this research. This 24 year period includes a historically significant bull market decade followed by an equally historically significant bear market decade.

Multiple vendor sources have been used in assembling the data so that single vendor bias can be eliminated. The construction of this point in time database has focused on collecting weekly consolidated data tables, collected every Friday from January 3, 1986 to the present, representing detailed point in time input to this study and cover the Value Line stocks on a weekly basis. Each stock record contains daily price and volume data, weekly analyst estimates and rankings, plus quarterly financial data as reported. The primary focus is on gross and net revenues.

Our historical database contains 1050 weeks of data between January 1990 and December 2009. In a full training and testing protocol there is a separate symbolic regression run for each of these 1050 weeks. Each SR run consists of predicting the *ftmEPS* for each of the Value Line stocks available in that week, using the 520 prior weeks as the training data set for that week. A sliding training/testing window will be constructed to follow a strict statistical out-of-sample testing protocol.

For each of the 1050 weeks, the 520 prior weeks training examples will be extracted from records in the historical trailing 10 years behind the selected record BUT *not including any data from the selected week or ahead in time*. The training dependent variable will be extracted from the historical data record exactly 52 weeks forward in time from the selected record BUT *not including any data from the selected week or ahead in time*. Thus, as a practical observation, the training will not include any records in the first 52 weeks prior to the selected record—*because that would require a training dependent variable which was not available at the time*.

For each of the 1050 weeks, the testing samples will be extracted from records in the historical trailing 10 years behind the selected record *including all data from the selected week BUT not ahead in time*. The testing dependent variable will be extracted from the historical data record exactly 52 weeks forward in time from the selected record.

M.F. Korns

Each experimental protocol will produce approximately ∼1500 linear regressions and 2 symbolic regression runs over an average of ∼780,000 (∼*1500 × 520*) records for each training run and for ∼1500 records for each testing run. Ten hours will be allocated for training. Of course separate R-Square statistics will be produced for each experimental protocol. We will examine the R-Square statistics for evidence favoring the addition of swarm intelligence over the base line and for evidence favoring one swarm intelligence technique over another.

Finally we will need to adapt our methodology to conform to the rigorous United States Securities and Exchange Commission oversight and regulations on investment managers. The SEC mandates that every investment firm have a compliance officer. For any automated forward earnings prediction algorithm, *which would be used as the basis for later stock recommendations to external clients or internal portfolio managers*, the computer software code used in each prediction, the historical data used in each prediction, and each historical prediction itself, must be filed with the compliance officer in such form and manner so as to allow a surprise SEC compliance audit to reproduce each individual forward prediction exactly as it was at the original time of publication to external clients or internal portfolio managers.

Of course this means that we must provide a copy of all code, all data, and each forward prediction for each stock in each of the 1050 weeks, to our compliance officer. Once management accepts our symbolic regression system, we will also have to provide a copy of all forward predictions on an ongoing basis to the compliance officer.

Furthermore there is an additional challenge in meeting these SEC compliance details. The normal manner of operating GP, and symbolic regression systems in academia will not be acceptable in a real world compliance environment. Normally, in academia, we recognize that symbolic regression is a heuristic search process and so we perform multiple SR runs, each starting with a different random number seed. We then report based on a statistical analysis of results across multiple runs. This approach produces *different results* each time the SR system is run. In a real world compliance environment such practice would subject us to serious monetary fines and also to jail time.

The SEC compliance requirements are far from arbitrary. Once management accepts such an SR system, the weekly automated predictions will influence the flow of millions and even billions of dollars into one stock or another and the historical back testing results will be used to sell prospective external clients and internal portfolio managers on using the system's predictions going forward.

First the authorities want to make sure that as time goes forward, *in the event that the predictions begin to perform poorly*, we will not simply rerun the original predictions again and again, with a different random number seed, until we obtain better historical performance and then substitute the new better performing historical performance results in our sales material.

Second the authorities want to make sure that, *in the event our firm should own many shares of the subsequently poorly performing stock of "ABC" Corp*, that we do not simply rerun the current week's predictions again and again, with a

different random number seed, until we obtain a higher ranking for "ABC" stock 221
thus improperly influencing our external clients and internal portfolio managers to 222
drive the price of "ABC" stock higher. 223

In order to meet SEC compliance regulations we have altered our symbolic 224
regression system, used in this chapter across all experiments, to use a pseudo 225
random number generator with a pre-specified starting seed. Multiple runs always 226
produce *exactly the same results*. 227

## 21.3  Investing Strategies 228

Value investing (Graham and Dodd 2008) has produced several of the wealthiest 229
investors in the world including Warren Buffet. Nevertheless, value investing has 230
a host of competing strategies including momentum (Bernstein 2001) and hedging 231
(Nicholas 2000). 232

One of the most difficult challenges in devising a securities investing strategy 233
is the a priori identification of pending regime changes. For instance, momentum 234
investing strategies were very profitable in the 1990s and not so profitable in the 235
2000s while value investing strategies were not so profitable in the 1990s but 236
turned profitable in the 2000s. Long Short hedging strategies were profitable in the 237
1990s and early 2000s but collapsed dramatically in the late 2007 thru 2008 period. 238
Knowing when to switch from Momentum to Value, Value to Hedging, and Hedging 239
back to Value was critical for making consistent above average profits during the 20 240
year period from 1990 thru 2009. 241

The challenge becomes even more difficult when one adds the numerous 242
technical and fundamental buy/sell triggers to currently popular active management 243
investing strategies. Bollinger Bands, MACD, Earning Surprises, etc. all have com- 244
plex and dramatic effects on the implementation of securities investing strategies, 245
and all are vulnerable to regime changes. The question arises, "*Is there a simple 246
securities investing strategy which is less vulnerable to regime changes than other* 247
*strategies*?". 248

An idealized value investing hypothesis is put forward: "*Given perfect foresight*, 249
*buying stocks with the best future earning yield (**Future12MoEPS/CurrentPrice**)* 250
*(ftmEP) and holding for 12 months will produce above average securities investing* 251
*returns*". 252

Of course the ideal hypothesis is *impossible to implement* because it requires 253
perfect foresight which is, in the absence of time travel, unobtainable. Nevertheless 254
the ideal hypothesis represents the theoretical upper limit on the profits realizable 255
from a strategy of buying future net revenue cheaply; yet, the theoretical profits 256
are so rich that one cannot help but ask the question, "*Are there revenue prediction* 257
*models which will allow one to capture some portion of the profits from the ideal* 258
*hypothesis*?". 259

The easiest revenue prediction model involves simply using the current year's 260
trailing 12 month revenue as a proxy for future revenue. 261

M.F. Korns

**Table 21.1** Returns for SP100 High ttmEP/ftmEP 100

| Year | SP100 stocks | 100ttmEP stocks | 100 ftmEP stocks | |
|------|------|------|------|------|
| 1990 | (6 %) | (17 %) | 3 % | t3.1 |
| 1991 | 24 % | 40 % | 111 % | t3.2 |
| 1992 | 3 % | 22 % | 56 % | t3.3 |
| 1993 | 8 % | 9 % | 46 % | t3.4 |
| 1994 | 0 % | 6 % | 18 % | t3.5 |
| 1995 | 36 % | 22 % | 49 % | t3.6 |
| 1996 | 23 % | 28 % | 38 % | t3.7 |
| 1997 | 28 % | 27 % | 51 % | t3.8 |
| 1998 | 32 % | 12 % | 12 % | t3.9 |
| 1999 | 31 % | 38 % | 22 % | t3.10 |
| 2000 | (13 %) | 14 % | 45 % | t3.11 |
| 2001 | (15 %) | 11 % | 56 % | t3.12 |
| 2002 | (24 %) | (15 %) | 8 % | t3.13 |
| 2003 | 24 % | 52 % | 67 % | t3.14 |
| 2004 | 4 % | 13 % | 45 % | t3.15 |
| 2005 | (1 %) | 17 % | 43 % | t3.16 |
| 2006 | 16 % | 7 % | 19 % | t3.17 |
| 2007 | 3 % | (5 %) | 20 % | t3.18 |
| 2008 | (37 %) | (28 %) | (17 %) | t3.19 |
| 2009 | 19 % | 43 % | 120 % | t3.20 |
| CAGR% | 6 % | 14 % | 37 % | t3.21 |
| Volatility | 20 % | 20 % | 30 % | t3.22 |
| CAGR% 1990s | 17 % | 18 % | 38 % | t3.23 |
| CAGR% 2000s | (4 %) | 8 % | 37 % | t3.24 |

*Note*: Per annum total returns for each year

The data supports the conclusion that even using this current revenue proxy model buying the top one hundred stocks with the highest (***current12MoEPS/ currentPrice***) (*ttmEP*) and holding for 1 year produces above average securities investing profits, *as least for the Value Line stocks*, as shown in Table 21.1.

Nevertheless, buying a stock with high EP, *but whose future 12 month earnings will plummet bringing on bankruptcy*, is an obviously poor choice. So why is high EP investing so successful given that future 12 month earnings can vary significantly? Placing current earnings yield investing in this context puts a new spin on this standard *value investing* measure. In this context we are saying that current earnings yield (also known as high EP investing) works precisely to the extent that *current earnings are a reasonable predictor of future earnings*! In situations where current earnings are NOT a good predictor of future earnings, then current earnings yield investing loses its efficacy.

This agrees with our common sense understanding. For instance, given two 275
stocks with the same high current earnings yield, where one will go bankrupt next 276
year and the other will double its earnings next year; we would prefer the stock 277
whose earnings will double. Implying that, *in the ideal*, current earnings are just a 278
data point. We want to buy *future earnings* cheap! 279

Precisely because the per annum returns from this current revenue prediction 280
model are far less than the returns achieved with perfect prescience, we must now 281
look for more accurate methods of net revenue prediction. 282

### 21.3.1  Estimating Forward 12 Month EPS 283

Each week we will perform ∼1500 linear regressions, one for each of the Value 284
Line stocks. The preapproved linear regressions are expressed by the following 285
Regression Query Language **RQL** (Korns 2013) expression: 286

$$\textbf{regress}\,(\textbf{x0}, \textbf{x1}, \textbf{x2}, \textbf{x3}, \textbf{x4}, \textbf{x5}, \textbf{x6})\;\textbf{where}\;\{\}$$

For each of the ∼1500 Value Line stocks in the current week, from each of the 287
520 trailing historical weeks for that stock (*see our methodology section above*) the 288
following seven input (*independent*) variables will be collected: 289

| | | | |
|---|---|---|---|
| 1. | **Estimated12MoEPS**(x0) | Wall Street analysts 12Mo forward EPS estimate | t6.1 |
| 2. | **Forward12MoEPS**(x1) | CurrentEPS + (CurrentEPS-Past1YrEPS) | t6.2 |
| 3. | **Projected12MoEPS**(x2) | CurrentEPS + ((CurrentEPS-Past1QtrEPS) * 4) | t6.3 |
| 4. | **EstimatedS12MoEPS**(x3) | (Wall Street analysts 12Mo forward SPS estimate) * CurrentMargin | t6.4 |
| 5. | **ForwardS12MoEPS**(x4) | (CurrentSPS + (CurrentSPS-Past1YrSPS)) * CurrentMargin | t6.5 |
| 6. | **ProjectedS12MoEPS**(x5) | (CurrentSPS + ((CurrentSPS-Past1QtrSPS) * 4)) * CurrentMargin | t6.6 |
| 7. | **WeeksSinceLastReport**(x6) | Absolute count of weeks since last quarterly report | t6.7 |

Each of the ∼1500 linear regressions produces an *ftmEPS* estimate for each of 290
the Value Line stocks for that week (**LRegress12MoEPS**). This regression output 291
is then used as an input to a single preapproved nonlinear weighted regression on 292
the following input variables: 293

The preapproved nonlinear weighted regression is expressed by the following 294
Regression Query Language **RQL** (Korns 2013) expression: 295

| 1. | **Estimated12MoEPS** ($x0$) | Wall Street analysts 12Mo forward EPS estimate | t9.1 |
|----|------|------|------|
| 2. | **Forward12MoEPS**($x1$) | CurrentEPS $+$ (CurrentEPS-Past1YrEPS) | t9.2 |
| 3. | **Projected12MoEPS**($x2$) | CurrentEPS $+$ ((CurrentEPS-Past1QtrEPS) * 4) | t9.3 |
| 4. | **EstimatedS12MoEPS**($x3$) | (Wall Street analysts 12Mo forward SPS estimate) * CurrentMargin | t9.4 |
| 5. | **ForwardS12MoEPS**($x4$) | (CurrentSPS $+$ (CurrentSPS-Past1YrSPS)) * CurrentMargin | t9.5 |
| 6. | **ProjectedS12MoEPS**($x5$) | (CurrentSPS $+$ ((CurrentSPS-Past1QtrSPS) * 4)) * CurrentMargin | t9.6 |
| 7. | **WeeksSinceLastReport**($x6$) | Absolute count of weeks since last quarterly report | t9.7 |
| 8. | **LRegress12MoEPS**($x7$) | Result of the linear regression for the stock in question | t9.8 |

$$\mathbf{model}\big(\mathbf{c0 * f0\,(x0, v0)}\,, \mathbf{c1 * f1\,(x1, v1)}\,, \mathbf{c2 * f2\,(x2, v2)}\,,$$

$$\mathbf{c3 * f3\,(x3, v3)}\,, \mathbf{c4 * f4\,(x4, v4)}\,, \mathbf{c5 * f5\,(x5, v5)}\,,$$

$$\mathbf{c6 * f6\,(x6, v6)}\,, \mathbf{c7 * f7\,(x7, v7)}\big)$$

$$\mathbf{where}\ \big\{\mathbf{op\,(\aleph, +, -, min, max)}$$

$$\mathbf{c0\,(0.0, 1.0)}\ \ \mathbf{c1\,(0.0, 1.0)}\ \ \mathbf{c2\,(0.0, 1.0)}\ \ \mathbf{c3\,(0.0, 1.0)}\ \ \mathbf{c4\,(0.0, 1.0)}$$

$$\mathbf{c5\,(0.0, 1.0)}\ \ \mathbf{c6\,(0.0, 1.0)}\ \ \mathbf{c7\,(0.0, 1.0)}\big\}$$

This nonlinear weighted regression will achieve regulatory and client preapproval because it is so intuitive and so easy to explain. Let us start with the simplest case where the functions (**f0** thru **f7**) are all noops $= \aleph$, then the final result will always be like the following example:

$$\mathbf{NLREstimated12MoEPS}\,(y) = \mathbf{.34 * x0} + \mathbf{.16 * x1} + \mathbf{.81 * x2} + \mathbf{.54 * x3}$$
$$+ \mathbf{.26 * x4} + \mathbf{.72 * x5} + \mathbf{.59 * x6} + \mathbf{.21 * x7}$$

We have eight inputs in the form of dollar values for next year's estimated EPS. Our model simply assigns a weight (*0.0 <= 1.0*) to each estimate—with the added benefit that, in the past 520 weeks for all ∼1500 Value Line stocks, these weights have been the most successful in predicting next year's EPS values for the Value Line stocks. Now moving on to the case where one of more of the functions (**f0** thru **f7**) are other than noops, then the final result will always be something like the following example:

$$\textbf{NLREstimated12MoEPS}\,(y) = \textbf{.34} * \textbf{x0} + \textbf{.16} * \textbf{x1} + \textbf{.81} * \textbf{max}\,(\textbf{x2}, \textbf{x0})$$
$$+ \textbf{.54} * \textbf{x3} + \textbf{.26} * \textbf{x4} + \textbf{.72} * \textbf{x5} + \textbf{.59} * \textbf{x6}$$
$$+ \textbf{.21} * \textbf{x7}$$

Again we have eight inputs in the form of dollar values for next year's estimated 307
EPS. Our model simply assigns a weight (*0.0<= 1.0*) to each possible simple 308
combination of those estimates—with the added benefit that, in the past 520 weeks, 309
these weights and these combinations have been the most successful in predicting 310
next year's EPS values for the SP100 stocks. 311

In all cases we are simply weighting simple estimates or simple combinations 312
of estimates with combinations that will never get unruly or out of hand and with 313
weights which will always remain safely between 0.0 and 1.0. For this intuitive 314
nonlinear weighted regression, Regulatory and client preapproval will be easy to 315
obtain. 316

## 21.3.2    *Estimating Forward 12Mo Total Return*                          317

A close examination of the Efficient Market Hypothesis (*EMH*) shows that the 318
expectation of rational investing decisions plays a significant role in the EMH 319
conclusions in favor of passive index investing. Therefore, in addition to attempting 320
to purchase cheap stocks (*via some estimate of future 12Mo earnings*), we would 321
also like to purchase stocks from sellers whose decisions may not be as rational as 322
the EMH might hope. 323

Normally each stock trades within its own average trading volume over the 324
course of weeks and months. This trading volume can be expressed as a percent 325
*WeeksVolume* = **(total number of shares traded today)/(total shares outstand-** 326
**ing)**. For any given stock there will be periods of calm when weekly trading 327
percent (*WeeksVolume*) is light compared to the its historical average, and periods 328
of frenzy when the weekly trading percent (*WeeksVolume*) is very high compared to 329
its historical average. Our assertion is that *when a trading frenzy is underway the* 330
*buyer and seller are less rational than on normal trading days*. 331

The following nine input factors (*each of which combines some measure of* 332
*trading frenzy or intrinsic value or both*) will be converted to z Scores (Anderson 333
et al. 2002) and are defined as follows: 334

First we see that **z Panic Level** is computed from the nonlinear regression future 335
12Mo EPS estimate divided by the week's closing price (*i.e. the estimated future* 336
*EPS yield*) times the percent of outstanding shares traded that week (*Weeks Volume*) 337
times the current week's trading percent as in comparison with the prior 52 weeks 338
trading percent (*Volume52WeekRange*). This input will be high when the estimated 339
future earnings yield is high (*the stock is cheap*), when a high percent of outstanding 340
shares traded this week (*Weeks Volume*), and when this week's trading volume is on 341
the high side compared to the previous 52 week trading history for this stock. This 342

M.F. Korns

| 1. | **zPanicLevel** (*x0*) | ((NLRFuture12MoEPS/WeeksClose) * WeeksVolume * Volume52WeekRange)) | t12.1 |
|---|---|---|---|
| 2. | **zPriceMomentum**(*x1*) | (Past52WeekReturn * WeeksVolume) | t12.2 |
| 3. | **zDollarVolume**(*x2*) | (WeeksVolume * Shares * WeeksClose) | t12.3 |
| 4. | **zFutureEPSYield**(*x3*) | (NLRFuture12MoEPS/WeeksClose) | t12.4 |
| 5. | **zSalesAttractiveness**(*x4*) | (Current12MoSPS/WeeksClose) * WeeksVolume | t12.5 |
| 6. | **zCurrentValuation**(*x5*) | (CurrentVPS/WeeksClose) | t12.6 |
| 7. | **zValuationAttractiveness**(*x6*) | (CurrentVPS/WeeksClose) * WeeksVolume | t12.7 |
| 8. | **zWallStreetRank**(*x7*) | Current Wall Street analysts ranking as a z Score | t12.8 |
| 9. | **zFinancialRank**(*x8*) | Current Wall Street financial ranking as a z Score | t12.9 |

is a stock selling on much higher volume than normal with a very cheap future earnings yield. We use this input as a measure of panic on the seller's side. Since each of these inputs are z Scores, a high value for this input indicates that this stock is in a greater trading frenzy relative to other stocks this week. 343 344 345 346

Second we see that **z Price Momentum** is computed from the stock's past 52 week total return (*Past52WeekReturn*) times the week's trading volume (*Weeks Volume*). This input will be high for stocks with strong momentum selling on high trading volume. This is a popular stock, and we use this input as a measure of price momentum on the buyer's side. Since each of these inputs are z Scores, a high value for this input indicates that this stock enjoys greater momentum relative to other stocks this week. 347 348 349 350 351 352 353

Third we see that **z Dollar Volume** is an estimate of the total dollar value of the shares traded this week. This is a popular stock, and we use this input as a measure of relative dollar flow through this stock as opposed to other stocks this week. Since each of these inputs are z Scores, a high value for this input indicates that more dollars are flowing through this stock than other stocks this week. 354 355 356 357 358

Fourth we see that **z Future EPS Yield** is a measure of how cheap the future 12Mo EPS estimate divided by the week's closing price (*i.e. the estimated future EPS yield*) is compared to other stocks this week. This input will be high when the estimated future earnings yield is high (*the stock is cheap*). Since each of these inputs are z Scores, a high value for this input indicates that this stock is cheaper relative to other stocks this week. 359 360 361 362 363 364

Fifth we see that **z Sales Attractiveness** is computed from the current 12Mo SPS divided by the week's closing price (*i.e. the current sales yield*) times the percent of outstanding shares traded that week (*Weeks Volume*). This input will be high when the current sales yield is high (*the stock is cheap*), and when a high percent of outstanding shares traded this week (*Weeks Volume*. This is a stock selling on high volume with a very cheap current sales yield. We use this input as a measure of attraction on the buyer's side. 365 366 367 368 369 370 371

Sixth we see that **z Current Valuation** is a measure of the current enterprise value divided by the week's closing price (*i.e. the current VPS yield*). This input will be high when the current VPS yield is high (*the stock is cheap*). Since each of these inputs are z Scores, a high value for this input indicates that this stock is cheaper relative to other stocks this week. 372 373 374 375 376

Seventh we see that **z Valuation Attractiveness** is a measure of the current 377
enterprise value divided by the week's closing price (*i.e. the current VPS yield*) 378
times this week's trading volume (*Weeks Volume*). This input will be high when the 379
current VPS yield is high (*the stock is cheap*), and trading volume is high. Since 380
each of these inputs are z Scores, a high value for this input indicates that this stock 381
is more attractive to buyers relative to other stocks this week. 382

Eighth we see that **z Wall Street Rank** is a measure of the current Wall Street 383
analysts' rank for this stock. This input will be high when the Wall Street analysts' 384
rank for this stock is high. Since each of these inputs are z Scores, a high value 385
for this input indicates that this stock enjoys a higher Wall Street analyst ranking 386
relative to other stocks this week. 387

Ninth we see that **z Financial Rank** is a measure of the current Wall Street 388
analysts' financial rank for this stock. This input will be high when the Wall Street 389
analysts' financial rank for this stock is high. Since each of these inputs are z Scores, 390
a high value for this input indicates that this stock enjoys a higher Wall Street analyst 391
financial ranking relative to other stocks this week. 392

The following single output factor (*what we train on*) will be converted to 393
sigmoid score (Kleinbaum et al. 2010; Anderson et al. 2002) and is defined as 394
follows: 395

| 1. | **sFuture12MoReturn** (*y*) | The actual Future 12Mo Total Return—as a sigmoid-score | t15.1 |
|----|------|------|------|

Obviously we are not trying to predict actual future 12 month total return so much 396
as we are trying to predict relative future 12 month total return. We don't really 397
need to know actual future total 12 month returns. We only need to select the 100 398
Value Line stocks with the highest *relative* estimated future total 12 month return. 399
This allows us the luxury of converting the output variable (**sFuture12MoReturn**) 400
to a sigmoid factor, which allows us to perform a nonlinear logistic regression 401
(Kleinbaum et al. 2010) of the following form. 402

$$\textbf{logit}\big(\textbf{f0}\,(\textbf{x0}, \textbf{v0})\,, \textbf{f1}\,(\textbf{x1}, \textbf{v1})\,, \textbf{f2}\,(\textbf{x2}, \textbf{v2})\,, \textbf{f3}\,(\textbf{x3}, \textbf{v3})\,, \textbf{f4}\,(\textbf{x4}, \textbf{v4})\,, \textbf{f5}\,(\textbf{x5}, \textbf{v5})\,,$$

$$\textbf{f6}\,(\textbf{x6}, \textbf{v6})\,, \textbf{f7}\,(\textbf{x7}, \textbf{v7})\,, \textbf{f8}\,(\textbf{x8}, \textbf{v8})\big)\ \textbf{where}\ \{\textbf{op}\,(\aleph, +, -, \textbf{min}, \textbf{max})\}$$

This simple and extremely intuitive nonlinear logistic regression will easily 403
win regulatory and client preapproval. First of all this nonlinear regression will 404
never produce unexpected or wild output. It will produce an orderly estimate for 405
(**sFuture12MoReturn**) which will always lie between 0.0 and 1.0 for each stock. In 406
essence, this nonlinear regression model will automatically rank each stock between 407
0.0 and 1.0 in terms of estimated future 12 month total return (*with 1.0 being the* 408
*most desirable and 0.0 being the least desirable*). Let us start with the simplest case 409
where the functions (**f0** thru **f8**) are all noops $= \aleph$, then the final result will always 410
be like the following example: 411

$$\textbf{sFuture12MoReturn}(y) = \textbf{sigmoid}\big(.34 * \textbf{x0} + .16 * \textbf{x1} + .81 * \textbf{x2}$$
$$+ .54 * \textbf{x3} + .26 * \textbf{x4} + .72 * \textbf{x5} + .59 * \textbf{x6}$$
$$+ .21 * \textbf{x7} + .91 * \textbf{x8}\big)$$

We have nine inputs in the form of z Scores for factors combining some measure of relative value and/or trading frenzy. Our model simply projects each weighted factor onto a relative ranking between 0.0 and 1.0—with the added benefit that, in the past 520 weeks, these weights have been the most successful in predicting next year's relative future 12Mo total return for the Value Line stocks.

Now moving on to the case where one of more of the functions (**f0** thru **f8**) are other than noops, then the final result will always be something like the following example:

$$\textbf{sFuture12MoReturn}(y) = \textbf{sigmoid}\big(.34 * \textbf{x0} + .16 * \textbf{x1} + .81 * \textbf{max}(\textbf{x2}, \textbf{x0})$$
$$+ .54 * \textbf{x3} + .26 * \textbf{x4} + .72 * \textbf{x5} + .59 * \textbf{x6}$$
$$+ .21 * \textbf{x7} + .91 * \textbf{x8}\big)$$

Again we have nine inputs in the form of z Scores for factors combining some measure of relative value and/or trading frenzy. Our model projects each weighted factor or simple combination of factors onto a relative ranking between 0.0 and 1.0—with the added benefit that, in the past 520 weeks, these weights and these combinations have been the most successful in predicting next year's relative future 12Mo total return for the Value Line stocks.

### 21.3.3 Historical Returns

Applying all of these tools, techniques, and factors to the task of creating our semi-passive VEP100 index fund, we perform our 1502 regression runs for the first week in each year from 1990 thru 2009. We select the 100 Value Line stocks with the highest **sFuture12MoReturn** values. And hold them for 1 year. We then compare the results to the SP100 passive index, buying the 100 stocks with the highest ttmEP, and buying the 100 stocks with the highest ftmEP and present the results in Table 21.2.

Our VEP100 semi-passive index produced a much higher compound annual growth rate (*CAGR%*) than the SP100 index and the 100 ttmEP method. However, it cannot compete with the ideal ftmEP method (*where one can see into the future*). Nevertheless the total return of our VEP100 semi-passive index is impressive and will definitely appeal to a wide range of high net worth clients.

**So have we beaten the Efficient Market Hypothesis**? With a little bit of humor I can answer with a definite **Yes** and **No**.

**Table 21.2**  Returns VEP 100

| Year | SP100 stocks | 100 ttmEP stocks | VEP100 index fund | 100 ftmEP stocks | |
|------|--------------|------------------|-------------------|------------------|--------|
| 1990 | (6 %) | (17 %) | (22 %) | 3 % | t18.1 |
| 1991 | 24 % | 40 % | 47 % | 111 % | t18.2 |
| 1992 | 3 % | 22 % | 33 % | 56 % | t18.3 |
| 1993 | 8 % | 9 % | 23 % | 46 % | t18.4 |
| 1994 | 0 % | 6 % | 0 % | 18 % | t18.5 |
| 1995 | 36 % | 22 % | 30 % | 49 % | t18.6 |
| 1996 | 23 % | 28 % | 24 % | 38 % | t18.7 |
| 1997 | 28 % | 27 % | 31 % | 51 % | t18.8 |
| 1998 | 32 % | 12 % | 0 % | 12 % | t18.9 |
| 1999 | 31 % | 38 % | 30 % | 22 % | t18.10 |
| 2000 | (13 %) | 14 % | 10 % | 45 % | t18.11 |
| 2001 | (15 %) | 11 % | 38 % | 56 % | t18.12 |
| 2002 | (24 %) | (15 %) | (6 %) | 8 % | t18.13 |
| 2003 | 24 % | 52 % | 62 % | 67 % | t18.14 |
| 2004 | 4 % | 13 % | 30 % | 45 % | t18.15 |
| 2005 | (1 %) | 17 % | 29 % | 43 % | t18.16 |
| 2006 | 16 % | 7 % | 8 % | 19 % | t18.17 |
| 2007 | 3 % | (5 %) | 13 % | 20 % | t18.18 |
| 2008 | (37 %) | (28 %) | (42 %) | (17 %) | t18.19 |
| 2009 | 19 % | 43 % | 131 % | 120 % | t18.20 |
| CAGR% | 6 % | 14 % | 17 % | 37 % | t18.21 |
| Volatility | 20 % | 20 % | 30 % | 30 % | t18.22 |
| CAGR% 1990s | 17 % | 18 % | 18 % | 38 % | t18.23 |
| CAGR% 2000s | (4 %) | 8 % | 20 % | 37 % | t18.24 |

*Note*: Per annum total returns for each year

**Yes**, because the VEP100 CAGR% of 17 % is a whopping 9 % per annum greater    441
than the SP100! This is a significant amount which will be of interest to a large class    442
of serious investors. Furthermore, the performance of the VEP100 is more consistent    443
across bull and bear decades with a CAGR % of 18 % in the bullish 1990s and a    444
CAGR% of 20 % in the bearish 2000s. Coupled with the transparent and intuitive    445
methodology of the VEP100, there is definite added value here.    446

**No**, because the EMH does not actually claim that one cannot make higher profits    447
than the indices. The EMH claims that one cannot increase returns without also    448
increasing volatility, and this is exactly what happens with the VEP100 semi-passive    449
index. Volatility increases from 20 % with the SP100 to 30 % with the VEP100. So    450
in an important way, the VEP100 is a classic confirmation of the Efficient Market    451
Hypothesis.    452

## 21.4  Summary

Advances in both the industrial strength and accuracy of Symbolic Regression packages can help overcome the resistance to SR in the investment finance industry. Management trust, regulatory approval, and client acceptance, are no longer the severe hurdles that they were in the past. Improvements in SR robustness, result invariance, demonstrable accuracy, and regression constraint languages, such as Regression Query Language **RQL** (Korns 2010, 2013, 2014), now support regulatory and client preapproval of important component SR processes.

In this research work, as series of cascade linear and nonlinear SR regressions are used to create a transparent semi-passive index fund with significantly higher returns, over the 1990–2009 two decade period, than its Standard &Poors 100 index benchmark. Because of its transparent and algorithmic nature, the new VEP100 semi-passive index fund could enjoy much lower costs than a standard active fund and yet enjoy attractive returns—costs similar in nature to the SP100 passive index fund.

Future research will focus on other semi-passive indices with performance tailored to various diverse client needs and requirements, and regulatory approval issues.

AQ4  ## References

Graham, Benjamin, and David Dodd. 2008. Securities Analysis. New York, New York, USA. McGraw-Hill.

Kennedy, J.; Eberhart, R. 1995. Particle Swarm Optimization. *Proceedings of IEEE International Conference on Neural Networks*. **IV**. pp. 1942–1948.

Korns, Michael F. 2007. Large-Scale, Time-Constrained Symbolic Regression-Classification. In Riolo, Rick, L, Soule, Terrance, and Wortzel, Bill, editors, Genetic Programming Theory and Practice V, New York, New York, USA. Springer, pp. 299–314.

Korns, Michael F., and Nunez, Loryfel, 2008. Profiling Symbolic Regression-Classification. In Riolo, Rick, L, Soule, Terrance, and Wortzel, Bill, editors, Genetic Programming Theory and Practice VI, New York, New York, USA. Springer, pp. 215–228.

Korns, Michael F., 2009. Symbolic Regression of Conditional Target Expressions. In Riolo, Rick, L, Soule, Terrance, and Wortzel, Bill, editors, Genetic Programming Theory and Practice VII, New York, New York, USA. Springer, pp. 211–228.

Korns, Michael F., 2010. Abstract Expression Grammar Symbolic Regression. In Riolo, Rick, L, Soule, Terrance, and Wortzel, Bill, editors, Genetic Programming Theory and Practice VIII, New York, New York, USA. Springer, pp. 109–128.

Price, Kenneth, Storn, Rainer, Lampinen, Jouni 2009. Differential Evolution: A Practical Approach to Global Optimization. New York, New York, USA. Springer.

Guido Smits, Ekaterina Vladislavleva, and Mark Kotanchek 2010, Scalable Symbolic Regression by Continuous Evolution with Very Small Populations, in Riolo, Rick, L, Soule, Terrance, and Wortzel, Bill, editors, *Genetic Programming Theory and Practice VIII*, New York, New York, USA. Springer, pp. 147–160.

Flor Castillo, Arthur Kordon, and Carlos Villa 2010, Genetic Programming Transforms in Linear Regression Situations, in Riolo, Rick, L, Soule, Terrance, and Wortzel, Bill, editors, *Genetic Programming Theory and Practice VIII*, New York, New York, USA. Springer, pp. 175–194.

Trent McConaghy, Pieter Palmers, Gao Peng, Michiel Steyaert, Goerges Gielen 2009, Variation-Aware Analog Structural Synthesis: A Computational Intelligence Approach. New York, New York, USA. Springer.

J.A., Nelder, and R. W. Wedderburn, 1972, *Journal of the Royal Statistical Society, Series A, General*, 135:370–384.

John R Koza 1992, Genetic Programming: On the Programming of Computers by Means of Natural Selection. Cambridge Massachusetts, The MIT Press.

Korns, Michael F., 2011a. Accuracy in Symbolic Regression. In Riolo, Rick, L, Soule, Terrance, and Wortzel, Bill, editors, Genetic Programming Theory and Practice IX, New York, New York, USA. Springer (*to be published in winter 2011*).

Pham, D., Ghanbarzadeh, A., Koc, E., Otri, S., Rahim, S., and Zaidi, M. 2005. "The Bees Algorithm". Technical Report Cardiff University.

Parpinelli, R. S., and Lopes, H. S., 2011. New inspirations in swarm intelligence: a survey. *Int Journal of Bio-inspired Computation*. **Vol 3**. Number 1.

Bernstein, J., 2001. Momentum Stock Selection: Using The Momentum Method for Maximum Profits. New York, New York, McGraw Hill

Nicholas, J., 2000. Market-Neutral Investing: Long/Short Hedge Fund Strategies. New York, New York, Bloomberg Press.

Poli, Riccardo, McPhee, Nicholas, Vanneshi, Leonardo, 2009. Analysis of the Effects of Elitism on Bloat in Linear and Tree-based Genetic Programming. In Riolo, Rick, L, Soule, Terrance, and Wortzel, Bill, editors, Genetic Programming Theory and Practice VI, New York, New York, USA. Springer, pp. 91–110.

Korns, Michael F, 2011b. Accuracy in Symbolic Regression. In Riolo, Rick, L, Soule, Terrance, and Wortzel, Bill, editors, Genetic Programming Theory and Practice IX, New York, New York, USA. Springer.

Korns, Michael F., 2012a. A Baseline Symbolic Regression Algorithm. In Soule, Terrance, and Wortzel, Bill, editors, Genetic Programming Theory and Practice X, New York, New York, USA. Springer.

Korns, Michael F., 2013. Extreme Accuracy in Symbolic Regression. In Soule, Terrance, and Wortzel, Bill, editors, Genetic Programming Theory and Practice XI, New York, New York, USA. Springer.

Korns, Michael F., 2014. Extremely Accurate Symbolic Regression for Large Feature Problems. In Soule, Terrance, and Wortzel, Bill, editors, Genetic Programming Theory and Practice XII, New York, New York, USA. Springer.

Korns, Michael F., 2012b. Predicting Corporate Forward 12 Month Earnings, 2012. Theory and New Applications of Swarm Intelligence, ISBN 978-953-51-0364-6, edited by Rafael Parpinelli and Heitor S. Lopes, InTech Academic Publishers.

Kleinbaum, David G., and Klein, Michael, 2010. Logistic Regression: A Self-Learning Text (Statistics for Biology and Health), ISBN 978–1441917416, New York, New York, USA. Springer.

Anderson, David R., Sweeney, Dennis J., and Williams, Thomas A, 2002. Essentials of Statistics for Business and Economics, ISBN 978–0324145809, Southwestern College Publishers.

AUTHOR QUERIES

AQ1.   Please note that the reference style has been changed from a Numbered style to a Name–Date style as per the style.
AQ2.   Please check the change made in the sentence "In addition to our own ARC system . . . " and correct if necessary.
AQ3.   Please provide department and institute/university details in the affiliation.
AQ4.   Refs. "Kennedy and Eberhart (1995), Korns (2007, 2009, 2011a, b), Korns et al. (2008), Price et al. (2009), Pham et al. (2005), Parpinelli and Lopes (2011), and Poli et al. (2009)" are not cited in the text. Please provide the citation or delete them from the list.