

Problems with scoring methods and ordinal scales in risk assessment

D. Hubbard
D. Evans

Risk assessment methods based on scoring methods that rate the severity of each risk factor on an ordinal scale are widely used and frequently perceived by users to have value. We argue that this perceived benefit is probably illusory in most cases. We begin by describing a number of common scoring methods currently used to assess risk in a variety of different domains. We then review the literature on the use of ordinal scales in risk analysis, the use of “verbal scales” for eliciting estimates of risks and probabilities, and the extensive research about peculiar human errors when assessing risks. We also supplement this overview with some data of our own. When these diverse kinds of evidence are combined, the case against scoring methods is difficult to deny. In addition to the evidence against the value of scoring methods, there is also a lack of good evidence in their favor. We conclude our overview by reviewing the reasons why risk assessment approaches should describe risk in terms of mathematical probabilities.

Introduction

Many methods for risk assessment involve the use of scoring methods in which the severity of each risk factor is rated on an ordinal scale. The resulting values are then combined by additive weighting or by multiplication to compute an aggregate measure of overall risk. In this paper, we provide an overview of the existing research and argue that, taken together, these diverse examples of evidence suggest that scoring methods are not useful tools for risk assessment. In arguing for this conclusion, we are not merely claiming that arithmetically combining ordinal scales is logically invalid. This is true regardless of risk management context. Rather, we claim that the widespread use of scoring methods in risk assessment is flawed for a wide variety of reasons. To a mathematically sophisticated reader, some of our points may appear unoriginal or obvious. This, however, would be to miss our primary point; the widespread use of scoring methods in real-world settings is still a serious problem that needs addressing, and these methods are beset by many flaws aside from the mathematical ones.

We start by describing a number of common scoring methods currently used to assess risk in a variety of different domains. We then identify and discuss four important

problems associated with the use of these scoring methods, which we may briefly state as follows. First, they do not usually take into account the findings of psychological research concerning the cognitive biases that impair most people’s ability to assess risk. Second, the verbal labels used in ordinal scales are interpreted extremely inconsistently among different users and by the same user. Third, many users treat these scales as if they are ratio scales, with the result that they draw invalid inferences. Fourth, simple scoring methods rarely consider correlations that would change the relative risks. Taken together, these four problems indicate that scoring methods are likely to be poor tools for risk assessment.

The only “evidence” in favor of using scoring methods in risk assessment consists of testimonials from users of such methods, claiming that they perceived great value in these methods. Published case studies even exist in the literature describing the application of such methods. However, subjective perceptions of benefit do not always correlate with objective measures of success. Such objective evidence in favor of scoring methods in risk assessment is lacking. When the lack of evidence in favor of such methods is combined with the powerful arguments against them, one is forced to conclude that scoring methods should not be used for risk assessment. We conclude this paper by arguing that risk

Digital Object Identifier: 10.1147/JRD.2010.2042914

© Copyright 2010 by International Business Machines Corporation. Copying in printed form for private use is permitted without payment of royalty provided that (1) each reproduction is done without alteration and (2) the Journal reference and IBM copyright notice are included on the first page. The title and abstract, but no other portions, of this paper may be copied by any means or distributed royalty free without further permission by computer-based and other information-service systems. Permission to republish any other portion of this paper must be obtained from the Editor.

0018-8646/10/\$5.00 © 2010 IBM

68 assessment should instead involve the use of mathematical
69 probabilities.

70 **Scoring methods for risk assessment:**
71 **A brief survey**

72 Many risk management methods promoted by management
73 consultants and international standard organizations involve
74 calculating a numerical value for risk based on simple point
75 scales (defined below) combined in some way. Such scales
76 are subjective and are usually based on some kind of ordinal
77 comparisons or classifications.

78 On an ordinal scale, factors such as likelihoods are
79 assigned numbers in such a way that the order of the numbers
80 reflects the order of the factors on an underlying attribute
81 scale. Two factors x and y with attribute values $a(x)$ and $a(y)$
82 are assigned numbers $m(x)$ and $m(y)$ such that $m(x) > m(y)$ if
83 $a(x) > a(y)$.

84 Two broad categories of subjective ordinal risk assessment
85 methods are widely used. We will refer to them as additive
86 and multiplicative. Additive scores are those in which several
87 independent factors are weighted and added together. For
88 example, when evaluating the credit risk of an international
89 customer, a manufacturer might identify factors for country
90 stability, order size, currency volatility, and so on. These
91 scores tend to be used to evaluate overall risks of projects,
92 investments, policies, etc. Multiplicative scores are those that
93 have just two or three factors that are multiplied together.
94 When two factors are used, they are usually *likelihood* and
95 *impact* (or probability and consequence), and they are
96 generally not individually weighted as in additive schemes.
97 These scores are often used to represent the evaluation of the
98 risk of individual events, such as theft of data by hackers
99 (i.e., malicious users) or power outage in the factory. Some
100 three-factor variations of multiplicative scores, particularly in
101 security-related areas, use *threat*, *vulnerability*, and *impact*.

102 In each of these approaches, the following
103 characterizations apply. First, some set of factors is
104 identified—possibly many factors for additive methods but
105 usually just two or three for multiplicative methods. Next,
106 each factor is rated according to some kind of point scheme.
107 A typical framework is an ordinal scale of 1–5, but other
108 scales are also used. If project duration is considered a risk
109 factor in an additive scoring method for information
110 technology (IT) project risks, then the duration in months
111 might be converted to discrete values on this five-point scale.
112 Sometimes, simple three-point *high*, *medium*, and *low* scales
113 are used. After rating each factor, the factors are then
114 combined. For additive scoring methods, weights are applied
115 to each factor, and then the weighted values are added. In
116 scoring methods that simply utilize likelihood and impact,
117 the weighted values are usually multiplied. Sometimes, the
118 ordinal scale is converted to yet another scale before the
119 weighted values are combined. For example, on a
120 multiplicative scale, a “low” impact score might be converted

to a 1 and a “high” impact converted to a 10. Finally, the
resulting value is compared with yet another arbitrary scale
(e.g., a score of 10–18 is “medium” risk, and a score of
greater than 25 is “extreme risk”). In the case of
multiplicative scores, the likelihood and impact might be
plotted on a so-called “risk matrix” in which different regions
are defined as a high/medium/low risk.

These sorts of weighted scoring methods are relatively easy
to create and teach. Consequently, the use of scoring
methods tends to spread quickly, and the methods have become
popular for several applications. Respected organizations
have designed such methods and represent them as “best
practice” for thousands of users. For example, the U.S. Army
has developed a weighted scoring-based method for evaluating
the risks of missions [1]. The U.S. Department of Health and
Human Services uses a weighted scoring method to determine
vaccine allocations in the event of a pandemic outbreak [2],
and the U.S. National Aeronautics and Space Administration
(NASA) uses a scoring method for assessments of risks in
manned and unmanned space missions [3, 4].

Additionally, a variety of influential standards that are now
considered as “best practice” in the information-technology
industry include proposed risk assessment methods based on
ordinal scores. For example, the National Institute of
Standards and Technology (NIST), which is a federal agency
of the U.S. tasked with setting measurement standards, has
developed the NIST 800-30 standard for assessing
information-technology security risks [5]. The Project
Management Institute (PMI), which describes itself as the
“world’s leading not-for-profit association for the project
management profession” on its website, has trained and
certified thousands of people in its methods that include a
scoring approach to evaluating project risk. PMI membership
includes many IT professionals as well as non-IT
professionals. In its Project Management Body of
Knowledge, PMI proposes its own risk scoring approach
based on ordinal scales [6]. Finally, the Information Systems
Audit and Control Association has certified thousands of
professionals in its training programs and has developed the
Val IT and Risk IT methods for evaluating the value and risk
of IT investments.

Of the cited examples above, all but two are additive
scores. The NASA and NIST scores are multiplicative.
The NIST score is an example of a scheme that uses
high/medium/low categories and then converts those
categories to another quantity (0.1, 0.5, and 1 for likelihood;
10, 50, and 100 for impact) before multiplying them
together.

Clearly, weighted scoring methods have been adopted by
influential organizations for assessing a variety of risky
decisions involving not only major investments and projects
but also public health and safety. If these methods are flawed,
as we argue in the following section, then the consequences
could have a wide effect.

175 **Four problems with scoring methods**
176 In this section, we argue that scoring methods are badly
177 flawed, and identify four associated problems. Before
178 examining each problem in turn, however, it is valuable to
179 ask why these methods have spread so widely if they are
180 really as problematic as we claim. One main reason may be
181 that these methods tend to be applied to decisions for which
182 the results are not immediately apparent, and the quality of
183 the decision is not obviously measurable. If an organization
184 uses a scoring method to rank the risks of some relatively
185 infrequent event (e.g., flu pandemics, engineering
186 catastrophes, or major IT security breaches), a very long time
187 may be required to determine whether the assessments
188 effectively track outcomes. In addition, this assumes that the
189 organization is systematically tracking the event and
190 statistically analyzing the results with respect to original
191 forecasts, which is not often the case. More often, managers
192 may be interpreting anecdotal outcomes in a way that
193 supports a particular assessment method. As a result, scoring
194 methods are rarely evaluated in a rigorous fashion, and they
195 persist despite their uselessness or even harmfulness.

196 From our review of the literature, we have identified four
197 important problems, which we mentioned in “Introduction,”
198 associated with the use of these scoring methods. We now
199 examine each of these four problems in turn.

200 **Partial review of literature on cognitive biases**

201 Several decades of empirical research in psychology have
202 revealed that people are ineffective at assessing risk.
203 Moreover, they depart from normative standards in certain
204 systematic ways. In other words, people are not only
205 irrational, they are also “predictably irrational” [7]. In this
206 section, we provide a brief and partial review of the literature
207 on cognitive bias and discuss some of the implications for the
208 use of scoring methods in risk assessment.

209 The systematic ways in which people make poor risk
210 assessments are often attributed to a well-documented set of
211 cognitive biases that systematically skew human judgment in
212 certain directions. For example, when estimating the frequency
213 of dangerous events, most people ignore available statistical
214 information about base rates and instead base their estimates
215 on their memories, which are biased toward vivid, unusual,
216 or emotionally charged examples. Researchers refer to this bias
217 as “the availability heuristic,” referring to the recall of
218 “available” memories. Most people also tend to assume that
219 individual random events are influenced by previous random
220 events (referred to as the *gambler’s fallacy*), and overestimate
221 the probability of good things happening to them compared
222 with other people (referred to as the *optimism bias*). Most
223 people also search for or interpret information in a way that
224 confirms their preconceptions (*confirmation bias*) and claim
225 more responsibility for successes than failures (*self-serving
226 bias*). Many other cognitive biases have been identified by
227 psychologists [8].

228 Other phenomena associated with the risk estimation
229 process may affect subjective responses. Anchoring, for
230 example, is the phenomenon of making estimates based on
231 adjustments to a previously conceptualized quantity. This
232 effect appears to occur even when the previous quantity is
233 unrelated to the new estimate. In one experiment, the
234 researchers conducted a survey of college students in which
235 one question asked for the last four digits of their social
236 security number. The very next question asked for an
237 estimate of the number of physicians in Manhattan. The
238 researchers found a correlation of 0.4 between these two
239 answers, indicating that estimates of physician numbers had
240 been influenced by the estimators’ social security numbers,
241 although these are clearly unrelated. The same researchers
242 found that questions that are logically identical, yet “framed”
243 differently (such as survival probabilities versus mortality
244 probabilities for medical procedures), can give rise to
245 significantly different responses. These effects are routinely
246 considered in the design of public opinion surveys but are not
247 considered in nearly all of the common risk assessment
248 methodologies.

249 Perhaps one of the most significant biases related to risk
250 assessments is that of overconfidence [9, 10]. Most people
251 consistently overestimate their certainty about a forecast. In
252 other words, if a manager believes that there is a 90% chance
253 some event will not occur, and this is tracked for a large
254 number of such events, then he will be correct much less
255 than 90% of the time. A review of several experiments shows
256 that when subjects believe that they have identified a range
257 that has a 90% chance of containing the correct value, their
258 track record shows that this range only contains the correct
259 value between 55% and 78% of the time. Since assessing risk
260 must involve assessing the likelihood of some event,
261 overconfidence means that subjective estimates of chance
262 will consistently be biased so that somewhat unlikely events
263 will be underestimated. In other words, if a person believes
264 that there is a 90% chance that a catastrophic event will not
265 happen, then the chance the event will occur may be much
266 higher than they believe. In the technical jargon, subjective
267 probabilities are often *miscalibrated*.

268 Another important phenomenon is not so much a
269 bias—which tends to be consistent, by definition—but an
270 inconsistency itself. A person asked to evaluate the same
271 item on different occasions tends to give different estimates,
272 although no change in information justifies the change in
273 judgment. Simple linear regressions of human judgments
274 based on various questions as independent variables in
275 a judgment model improve upon unaided human judgment
276 [11, 12]. In other words, when people use such linear models,
277 they are more consistent than when they rely entirely on
278 intuition, and this alone improves performance. One
279 important reason for this is that the linear model removes
280 human inconsistency. Studies of this “intrapersonal
281 reliability” show it to be fairly low [13]. For example, a study

282 of X-ray specialists showed that the same set of cases
283 examined a week apart had a correlation of the prognoses of
284 just 0.6, indicating that specialists may make quite different
285 prognoses when presented with the same cases [14]. The
286 same specialists were giving prognoses that were different for
287 no reason other than random variation.

288 These biases and related challenges are persistent and are
289 not overcome without specific and extensive training,
290 controls, and adjustments. Some progress has been made in
291 devising mental techniques and tools that people can use to
292 overcome or at least reduce bias and thereby make more
293 rational decisions. These debiasing strategies include
294 meta-cognition, instruction in rational choice theory and
295 probability, decision-support systems, and cognitive forcing
296 [15]. Inconsistency seems to be partially offset by repeatedly
297 asking a number of similar questions and statistically
298 assessing the variance. The effects of response item order
299 (due to anchoring) can be offset by asking different persons
300 the same questions in different orders.

301 There is much more literature on heuristics and biases than
302 would fit within the scope of this paper. However, despite the
303 significance of the research on cognitive biases and its
304 relevance to risk assessment, none of the common scoring
305 methods account for such bias. Because assessing likelihood
306 is so central to risk assessments, *overconfidence* alone would
307 be a debilitating shortcoming for any risk assessment if left
308 unaddressed. The effects of anchoring, framing, and
309 inconsistency would only further compound this challenge.

310 **Existing research regarding the variability of** 311 **verbal labels**

312 A characteristic of virtually all scoring methods is the use of
313 some verbal labels for eliciting judgments. Likelihood, for
314 example, might be put into categories of “very likely,”
315 “likely,” “as likely as not,” “unlikely,” and “very unlikely.”
316 Additionally, as in the case of the NIST 800-30 security
317 model, an even simpler “high,” “medium,” or “low”
318 differentiation may be used. Similarly, such verbal
319 expressions are used to categorize impact (e.g., “moderate”
320 or “extremely severe”). Such kinds of verbal methods are
321 used because it is believed that users of these methods will
322 have better comprehension and, therefore, more reliable use
323 of the method if simple labels are used. It is also assumed that
324 if the labels are further defined in great detail, then the scales
325 will be used consistently. Unfortunately, neither of these
326 assumptions is valid when scrutinized in an empirical
327 fashion. In this section, we review some of the existing
328 evidence that undermines these assumptions.

329 For example, the Intergovernment Panel on Climate
330 Change report made use of verbal expressions of likelihood
331 instead of quantitatively expressed probabilities. A table was
332 provided to define what each of the labels meant in
333 quantitative terms. The table would define “very likely” as
334 meaning “greater than 90%” and “unlikely” as meaning “less

than 33%.” Research showed that when given specific uses
in context (e.g., “It is very likely that hot extremes, heat
waves, and heavy precipitation events will continue to
become more frequent”) and even when provided with
detailed definitions of the probability labels, subjects who
read the report interpreted the labels to mean a wide variety
of possible values even so far as to violate the defined
rule [16]. In other words, even when users were given tables
that explained what range of values were valid for specific
labels, users still interpreted the labels differently—sometimes
even outside of the specified ranges. It was even possible
for the label “unlikely” to be interpreted as meaning as much
as a 66% probability.

In the same manner, the NIST 800-30 guideline for
assessing IT security risks defines high, medium, and low
for both likelihood and impact. The definition of “low
likelihood” states that “the threat-source lacks motivation or
capability, or controls are in place to prevent, or at least
significantly impede, the vulnerability from being exercised”
and goes on to state that this will be converted to a likelihood
of “0.1.”

When different individuals interpret the same labels to
mean very different things, there can arise an “illusion of
communication” [16]. Subjects may describe the probability
of a given event with the same verbal label and conclude
on this basis that they agree; however, since they may
implicitly attach different probability ranges to the verbal
label, their agreement may be illusory. To further complicate
matters, the same individual may attach a different
probability range to the same label in different contexts. The
additional variance caused by the inconsistent interpretation
of these labels—even when care is taken to provide detailed
definitions—is another source of error in most scoring
models.

369 **Invalid inferences**

370 Arbitrary features of the scoring scheme itself often have a
371 larger impact on results than the users might be aware of.
372 In the first place, people who use them treat the values arrived
373 at by means of these methods as if they were values on a ratio
374 scale, when in fact they are, of course, values on an ordinal
375 scale. This lack of care leads users to draw invalid inferences
376 that are always misleading and sometimes harmful.

377 In 1946, Stevens introduced the idea of different “scales”
378 of measurement [17]. The scales describe the nature of
379 information contained within the numbers assigned to
380 attributes. Among the various levels of measurement
381 identified by Stevens were ordinal and ratio scales. Note that
382 in this paper we limit our discussion to the choice between
383 ordinal and ratio scales. Additionally, we do not believe that
384 the nominal and interval scales that Stevens also identified
385 are relevant to the debate on methods for risk assessment.
386 Unlike ordinal scales (described earlier), ratio scales assign
387 numbers to attributes so that differences and ratios between

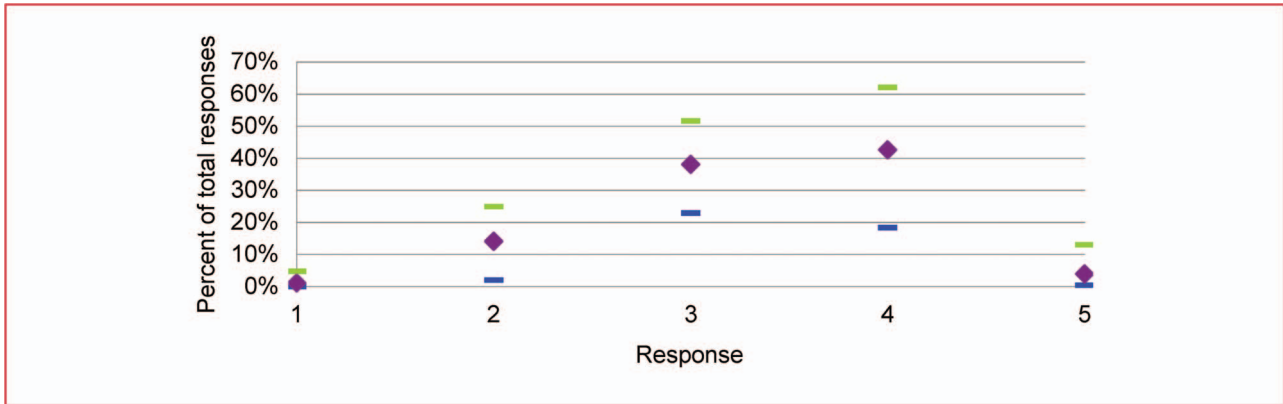


Figure 1

Distribution of item responses.

the numbers are meaningful; that is, such a scale has a natural zero value and is unique up to the choice of the scale unit. In other words, the scale is nonarbitrary.

We believe that many people tend to treat measurements as if they were made on a ratio scale unless explicitly warned to do otherwise. A certain kind of mental discipline and mathematical sophistication is needed to remember that one is using an ordinal scale, and we suspect that this discipline and sophistication is often lacking in those who use scoring methods in risk assessments. This can lead users to draw invalid inferences that may sometimes be harmful.

For example, likelihood may originally be assessed on a five-point verbal scale and later converted into a numerical representation for ease of processing. If “not very likely” is converted into a value of 2, and “very likely” is converted into a value of 4, then someone who treats this ordinal scale as if it were a ratio scale might infer that the likelihood of the latter category is exactly double the likelihood of the former. This inference may well be invalid, since ordinal scales involve an indeterminate distance metric.

“Range compression” is another problem that has been pointed out regarding the use of such scales [18]. Since ordinal scores tend to use scales with a small number of discrete values, there can be a significant loss of resolution in the assessment. This loss of resolution is more apparent when, as is done in some scoring models, known quantities are converted to an ordinal scale. For example, an additive scoring model that assesses the risk of IT projects may convert “planned duration” in months to a five-point ordinal scale such that less than three months is a “1,” three to six months is a “2,” and so on. One scoring method that has been proposed for ranking IT project portfolios converts return on investment to an ordinal scale so that each scale increment is an extremely wide range of values (e.g., 1%–299% is converted to a score of “1”) [19].

If the original quantity were uniformly distributed among the discrete score values, then the additional error added by this rounding off (i.e., conversion of wide ranges of values to a few scores) would be half of the range of values in one ordinal increment. Note that this rounding error is exacerbated by two factors. First, many of the scales defined are not linear (a “1” for project duration may mean one to three months, while a “4” may mean one to two years). Second, responses in the purely subjective scales are often highly clustered. One of the authors collected the data for scoring responses from the scoring schemes used in five different organizations. All were used for IT project risk assessment, and all were based on five-point scales. In each case, 3–12 individuals had to assess up to 40 factors for each member in a list of separate IT risks or IT projects for a total of more than 2,000 individual item responses. The distribution of responses for all five methods is shown in **Figure 1** (the bars indicate the minimum, average, and maximum among the seven methods). As may be observed, just two of the values on the scale account for the majority of choices. In effect, the scale is used more like a two- or three-point scale, further reducing the resolution and increasing the rounding error. This is a cursory examination of a few real-world examples, but it does make the nature of this problem clear.

Another assumption is that functions based on simple additive or multiplicative operations are an adequate approximation of a function that might realistically represent risks. This assumption may have an analogy with the fanciful example of aircraft engineers attempting to accurately compute the range of an aircraft by making use of weighting and adding factors associated with wingspan and fuel tank size instead of using the relevant differential equations. Actuaries and risk analysts who use quantitative simulations routinely need to use higher-order nonlinear equations that

458 are either derived from fundamental axioms of probability
459 theory or empirically validated.

460 **Invisible correlations**

461 More elaborate quantitative models of risks must generally
462 account for a variety of possible correlations among the
463 modeled variables. A portfolio that has 50 investments is not
464 well diversified if there are high correlations among all of
465 them (e.g., if they all involve parts for oil rigs, and their
466 business is a function of oil prices). Excluding such
467 correlations from the model would cause a systematic
468 underestimation of risks, and it would be considered a
469 significant oversight for such models to exclude correlations.
470 However, such issues are almost never considered in any of
471 the ordinal scoring methods.

472 Furthermore, it is common in quantitative models of the
473 risks of systems (e.g., nuclear power plants) to model
474 “cascade failures” and “common mode failures.” These are
475 single failures that simultaneously cause several other
476 failures and failures that cause a chain reaction sequence of
477 other failures, respectively. As with correlations, leaving
478 such relationships out of a model would be considered a
479 significant shortcoming that, again, results in the systematic
480 and dramatic underestimation of risks.

481 A multiplicative model, for example, will assess the
482 likelihood and impact of several different risk events as if they
483 were entirely independent. When the scores are aggregated,
484 three different events may be considered to be “medium” or
485 “low” risks. However, if the events are correlated and have an
486 impact on one another, then three low-to-medium risks can
487 together produce one very high risk.

488 **Describing risk in terms of probabilities**

489 We may summarize the four problems we identified in the
490 previous section in terms of four common fallacies that have
491 arisen about scoring methods. The first fallacy is the idea that
492 without further aid or adjustment, subjective assessments
493 of experts are an adequate assessment of risks. This is not
494 true; the research in human judgment and decision making
495 cited in the section “A partial review of literature on
496 cognitive biases” shows that various known errors and biases
497 are highly significant in the subjective assessment of risks
498 even for—in fact particularly for—experts.

499 The second fallacy is the idea that if a scale is very
500 “rigorously” defined, then it will be used in a predictable
501 consistent manner. This is not true either; the research cited in
502 the section “Existing research regarding the variability of
503 verbal labels” shows that the verbal scales are used in a highly
504 inconsistent manner, even when users are given a great deal of
505 descriptive detail about the definitions of the values.

506 The third fallacy is the idea that the exact scales and
507 scoring calculations chosen are arbitrary and will not
508 ultimately have as much bearing on the outcomes as the input
509 of the experts. This is patently false; the arbitrary choice of

whether a scale is a three- or five-point scale, together with
the choice of how the scales are to be combined, can affect
the results or produce nonsensical results. Furthermore, the
arbitrary application of operations like addition and
multiplication to these ordinal scales adds even more error.

The fourth and final fallacy is the idea that correlations and
other interactions among events are not necessary to model a
rough approximation. This is false: in quantitative risk
models, it is known that excluding correlations and
relationships such as cascade failures or common mode
failures can cause a model to greatly underestimate risks.
There is no reason to believe this issue is somehow alleviated
by using ordinal scales.

All of the previously identified scoring methods are
associated with these fallacies. However, this is also true for
the more “rigorous” approaches, such as the analytic
hierarchy process (AHP), that mathematically attempt to
analyze verbal preferences. Since the AHP can be used to
estimate the coefficients of ordinal scales, proponents of AHP
like to point out that such scales will then directly yield a
cardinal product that is potentially accurate and relevant.
However, this does not imply that AHP thereby avoids the
problems we have identified with the use of ordinal scales in
risk assessment, since these mathematical adjustments are
only valid for converting ordinal utility measures to cardinal
utility values. However, in risk assessment, subjective utility
is only part of the concern; without objectively valid
(i.e., well-calibrated) probability assessments and cost
estimates, even the most finely tuned subjective utility
values will be useless. The theory behind converting ordinal
utility into cardinal utility (due to von Neumann and
Morgenstern [20]) is quite a different topic than the degree
to which the predictive power of subjectively chosen scores
is empirically validated.

Likewise, the theoretical soundness of AHP (which has
been disputed [21–27]) is a separate issue from the empirical
description of whether the potential accuracy of AHP is ever
realized in practice. Regardless of how much personal
satisfaction that people may derive from using AHP, its
objective value as a risk assessment tool depends on showing
that experts using AHP outperform experts using their own
unaided intuition in a controlled forecasting experiment.
There is, however, no such evidence. Instead, the research
concerning the validity of AHP, as with other scoring
models, has focused on other measures, such as how well the
method predicts the *stated preferences* of users, and not how
well it forecasts real-world events, such as costs, project
failures, and industrial accidents [28]. Furthermore, nothing
in AHP solves the problems described in the section
“Existing research regarding the variability of verbal labels.”
The “garbage in/garbage out” problem (i.e., erroneous input
produces erroneous output) exists for AHP if eliciting
probabilities and risks by scales (in this case, pairwise
comparisons) is flawed.

564 Again, care must be taken not to confuse the perception of
565 a benefit with actual improvements in the accuracy of
566 assessments. Studies have already shown that gathering more
567 information and interacting with other individuals before
568 making a decision can improve subjective confidence in the
569 decision without improving its objective quality; indeed, it
570 can even decrease the objective quality [29, 30]. One
571 empirical study of decision quality does, in fact, show that
572 the use of a well-known AHP tool improved confidence
573 without an improvement—or perhaps even a degradation—in
574 decision quality [31].

575 To avoid all of these problems, risk assessment should use
576 methods that meet the following three criteria. First and
577 foremost, useful risk assessment methods must use explicit
578 probabilities and magnitudes of losses expressed
579 quantitatively instead of using surrogate verbal or ordinal
580 scales. In other words, instead of stating that likelihood and
581 impact are “high” and “medium,” we would state that there is a
582 “10% chance of a loss of inventory worth \$2 million
583 to \$4.5 million.” One objection to this approach is that the data
584 may not be available. However, the previously cited research
585 shows that experts can be trained to subjectively estimate
586 the values. Note that ordinal scales, which also depend on
587 subjective assessments, do not in any way alleviate the
588 problem associated with a lack of data; on the contrary, they
589 merely hide it. With quantitative probability estimates,
590 however, uncertainty is explicitly incorporated. Furthermore,
591 these estimates can be refined (given the existence of the
592 extensive research about human performance in assessing
593 subjective probabilities) and tracked against actual outcomes.

594 The second criterion that good risk assessment methods
595 should satisfy is that they should use Monte Carlo
596 simulations to explicitly model the relationships of
597 components in a system and their correlations. In other
598 words, instead of just adding or multiplying abstract factors,
599 realistically descriptive functions are defined. If the risks of a
600 major engineering project missing a deadline were being
601 assessed, then the function may involve the rules of a detailed
602 work plan with a range of durations defined for inclement
603 weather or accidents. This exhibits no more of a “garbage
604 in/garbage out” problem than any of the ordinal scoring
605 methods, but it does allow for explicit known mathematical
606 relationships (such as correlations that may exist among
607 the components of a chemical processing plant, or the effect
608 of demographic trends on a new mall) for which ordinal
609 scales hardly provide even a gross approximation.

610 Finally, the third criterion for good risk assessment methods
611 is that they should allow for research in human judgment and
612 decision making to be applied in a corrective fashion. Many of
613 the probabilities required for such quantitative models may
614 still rely on the subjective judgments of experts, but methods
615 can be employed to correct for biases, such as overconfidence.
616 For example, the U.K. government explicitly acknowledges
617 that optimism bias is a problem in planning and budgeting and

618 has developed measures for how to deal with optimism bias in
619 the government. The U.K. Department for Transport requires
620 project planners to use so-called “optimism bias uplifts” for
621 large transport projects to arrive at accurate budgets for
622 planned ventures [32].

623 The first two criteria involve the use of methods already
624 widely used in actuarial science, probability theory, statistics,
625 and decision sciences, where risk is thought of as the
626 probability of a loss of some magnitude. Industries associated
627 with nuclear power, oil exploration, weather models, and
628 actuarial science already make extensive use of such
629 methods. One key advantage of such methods is that they
630 produce results that can be compared with observation.
631 While it ~~is may not be~~ possible to validate whether an
632 assessment of a “medium” probability was realistic (since it
633 is ambiguous), probabilities can at least be compared with
634 observed frequencies. Furthermore, if probability functions
635 are known for input variables, probability theory and
636 simulations can be used to assess the risks of events that have
637 not yet occurred (e.g., the 1-in-500-year events routinely
638 assessed in studies of nuclear power safety).

639 We should make note of one clarification. In these
640 quantitative solutions to risk assessments, sometimes
641 probability and magnitude are multiplied together. This may,
642 at first, seem like a multiplicative version of the ordinal
643 scales, but in this case, neither dimension is ordinal.
644 Probabilities are measured quantitatively as real number
645 values between 0 and 1, inclusive, which is a ratio scale.
646 Magnitude of loss, likewise, is a ratio scale, not an ordinal
647 scale. Loss could be expressed as a monetary value or, as in
648 the case of public health and safety, numbers of injuries or
649 deaths.

650 **Validated methods for eliciting** 651 **subjective probabilities**

652 It is a common belief among users of ordinal scoring
653 methods that their situations are uniquely complex or lack the
654 large volumes of data they imagine are at the disposal of
655 analysts in other environments. They apparently conclude,
656 then, that the use of subjective ordinal scores is their only
657 recourse. However, methods exist that make the subjective
658 assessments of probabilities practical while controlling for
659 many of the sources of error mentioned previously. Three
660 such methods appear to significantly improve the ability of
661 experts to provide subjective estimates.

662 First, regular feedback can be provided to experts with
663 incentives for improved performance. Weather forecasters
664 using incentive systems based on the *Brier score* have
665 developed a highly calibrated sense of subjective
666 probabilities [33]. That is, when they say there is a 95%
667 chance of precipitation, it rains 95% of the time. While
668 weather forecasters have the advantage of empirical data
669 and detailed weather models, the forecasts often require
670 subjective assessments (and the original research in this area,

671 in the late 1970s, predates the more advanced weather
672 models used today).

673 Second, experts can be provided with so-called *calibration*
674 *training*. There is evidence of some success in using various
675 training techniques to calibrate experts to aid in assessing
676 subjective odds. One of the authors has found that a half-day
677 of training significantly reduces both overconfidence and
678 underconfidence [34]. Other researchers have found that
679 similar simple techniques have a significant effect on the
680 calibration of probabilities [35]. Training methods typically
681 involve repeated tests with trivia questions and feedback on
682 results. Specific techniques involve treating the probability
683 assessment in such a way that it is compared with a
684 corresponding bet involving known chances. One such
685 method is sometimes called the *equivalent urn* method [36].
686 If an expert believes an event has a 20% chance of
687 occurrence, then a bet on that event should be considered
688 equivalent to a bet based on whether a red marble will be
689 randomly selected from an urn filled with marbles of which
690 20% are red and 80% are green.

691 Finally, experts can participate in prediction markets.
692 When a number of people trade contingent assets whose
693 value is tied to a future event such as “Barack Obama will
694 win a second term of office,” the current price at which these
695 assets change hands provides a good indication of the
696 probability of the relevant event. This approach seems to be
697 reasonably well calibrated even if the market involves only
698 “play” money and a scenario in which players compete for
699 points [37].

700 As the previous citations illustrate, these methods of
701 eliciting subjective probabilities are well-researched
702 alternatives to ordinal scales. Such methods can and have
703 been used in any environment where historical data are not
704 considered sufficient for quantitative analysis and subjective
705 estimates are required. These methods avoid the ambiguity
706 error introduced by using verbal labels and can directly be
707 used in probabilistic simulations.

708 However, even users of more sophisticated simulations
709 might not be aware of the benefits of these methods. One of
710 the authors has conducted a survey of users of common
711 personal computer Excel** spreadsheet-based Monte Carlo
712 simulation tools (such as *@Risk* and *Crystal Ball*) and finds
713 that while subjective probability estimates are common in
714 Monte Carlo simulations, calibration training is rarely used.
715 From a total of 72 models created by 34 individual analysts,
716 most involved a significant number of subjective probability
717 estimates, but none used calibration training.

718 Currently, one of the authors has provided probability
719 calibration training to a total of 190 individuals, all of whom
720 completed training in which they were asked to provide a
721 probability they will be correct in the answers for a series of
722 true/false trivia questions. Of those, 166 also completed
723 training in which subjects were asked to provide subjective
724 90% confidence intervals for a variety of trivia questions.

This is similar to the methods used in previously cited
research on overconfidence. When these individuals were
initially tested, on average, they correctly answered only 47%
of the answers within their stated 90% confidence intervals.
After training and three to five rounds of additional tests, they
on average correctly answered 80% of the questions within
their 90% confidence intervals, and 66% of the subjects
were within the statistically allowable error for perfect
calibration. The calibration of binary probabilities showed that
subjects who assessed probabilities produced an expected
percentage correct to within 5% of their actual percentage
correct (i.e., they would be correct about 80% of the time when
they said they were 80% confident). Initially, they were
overconfident by an average of 14% (i.e., when they said
they were 90% confident they were correct only 76% of
the time). These results are more encouraging than some of the
attempts at calibration in the previous literature. Differences
may be due to the fact that this training employed a
combination of several calibration techniques, whereas in
the previous literature subjects were typically tested one
method at a time. In addition, the group of subjects for this
training, unlike in some of the previous calibration research,
consisted entirely of professionals who are routinely
confronted with real-world estimation problems.

Evidence of improvements on ordinal scores

There is little detailed analysis on the performance of ordinal
scales since the results are rarely tracked (and this may very
well explain why ordinal scales persist). Even most users of
detailed Monte Carlo simulations rarely keep track of
forecasting performance. However, we do know of at least
one interesting example for which a fairly large sample of
situations have been assessed with both a multiplicative
ordinal “risk map” and quantitative models with Monte Carlo
simulations, and which therefore permits a comparison of the
two types of approaches to risk assessment. This example
concerns NASA, which requires engineers and mission
managers to develop a “5 × 5 risk matrix” for each new
mission or space vehicle to identify key risks that may result
in partial or complete mission failure [4]. At the same time,
the NASA Headquarters Cost Analysis Division (CAD) has
been conducting a Monte Carlo-based analysis of the same
missions. The Monte Carlo simulations address cost and
schedule risks as well as the risk of various mission failures.
All of these missions have also been part of a historical
analysis in which “best fit” regression models compare
indicators of mission complexity, budget constraints, and
time constraints to actual mission failures.

The Monte Carlo simulation produced by the CAD uses
data that are based, in part, on a logistic regression of
previous mission failure rates with respect to a “complexity
index” consisting of 37 factors [38]. Some subjective
estimates of probabilities were also used when historical
data were insufficient. The simulations outperformed the

778 5 × 5 matrices in predicting partial and complete mission
779 failures. While several observed causes of mission failures
780 were identified by the group that conducted the Monte Carlo
781 simulations, *not one* actual mission failure was ever
782 anticipated by the groups that created 5 × 5 matrices.

783 It is worth noting that the Monte Carlo models also
784 consistently outperformed traditional deterministic estimates
785 of costs and duration. Just over half of the simulations
786 produced cost and schedule results that were within 10% of
787 actual observations, whereas the deterministic methods
788 consistently underestimated costs and schedule delays by
789 25%–55%.

790 The analysts of the CAD also note an important reaction
791 by mission managers and engineers regarding the use of
792 models that depend in part on quantitative historical data.
793 There is a tendency to think of each mission as unique, so
794 that historical data cannot be used in a meaningful way to
795 assess the risks of the new vehicle or mission. In our
796 experience, we have also found this opinion to be not
797 uncommon among managers in information technology,
798 business operations, and a variety of other fields. However,
799 the data clearly show that, in the case of NASA, the
800 quantitative models that were based on historical data
801 analysis and quantitative subjective estimates outperform the
802 qualitative judgments of experts.

803 While previously published research comparing the
804 performance of various risk analysis methods involving
805 real-world problems is scant, what does exist seems to
806 support the experience at NASA. One study in the oil
807 exploration industry showed that the sophistication of
808 the quantitative risk assessments correlated with various
809 measures of financial performance of oil exploration
810 companies [39]. Furthermore, the effect on financial
811 performance was observed to occur *just after* the more
812 advanced methods were adopted [40]. These researchers
813 equated “more advanced” with the use of quantitative risk
814 estimates, decision trees, Monte Carlo simulations, options
815 theory, and portfolio theory.

816 Conclusion

817 In this overview, we have argued that scoring methods are
818 poor tools for risk assessment. We first described a number of
819 common scoring methods currently used to assess risk in a
820 variety of different domains. We then identified and
821 discussed four important problems associated with the use of
822 these scoring methods. We concluded by arguing that risk
823 assessment should describe risk in terms of mathematical
824 probabilities.

825 The problem discussed in this paper is serious. The fact
826 that simple scoring methods are easy to use, combined with
827 the difficulty and time delay in tracking results with respect to
828 reality, means that the proliferation of such methods may
829 well be due entirely to their perceived benefits and yet have
830 no objective value.

A practical alternative to scoring methods exists, namely,
the use of quantitative probability estimates. Methods
involving such probability estimates have been validated by
observation and, at a minimum, control for many of the
sources of error. We therefore conclude that caution should
be exercised before continually using scoring methods to
make critical risk assessments. The size and scope of many of
the investment decisions and public health and safety
decisions easily justify a more rigorous analysis than is
permitted by the use of scoring methods involving ordinal
scales. At a minimum, rigorous empirical tests of the
predictive validity of a scoring method should be required
before they are applied to such critical risk analysis problems.

** Trademark, service mark, or registered trademark of Microsoft Corporation in the U.S., other countries, or both.

References

1. J. D. Van Vactor, “Risk mitigation through a composite risk management process: The U.S. army risk assessment,” *Org. Dev. J.*, vol. 27, no. 4, pp. 84–97, 2009.
2. United States Department of Health and Human Services, Draft Guidance on Allocating and Targeting Pandemic Influenza Vaccine, 2007. [Online]. Available: <http://www.pandemicflu.gov/vaccine/prioritization.html>
3. J. V. Turner, “Risk management on the space shuttle upgrades development program,” in *Proc. Joint ESA-NASA Space-Flight Safety Conf.*, B. Battrick and C. Preyssi, Eds., 2002, European Space Agency, ESA SP-486.
4. D. Hubbard, *The Failure of Risk Management: Why It's Broken and How to Fix it*. Hoboken, NJ: Wiley, 2009, pp. 237–238.
5. *Risk Management Guide for Information Technology Systems, NIST 800-30*, Nat. Inst. Stand. Technol., Gaithersburg, MD, 2002. [Online]. Available: <http://csrc.nist.gov/publications/nistpubs/800-30/sp800-30.pdf>
6. *A Guide to the Project Management Body of Knowledge*, Project Manage. Inst., Newtown Square, PA, 2008.
7. D. Ariely, *Predictably Irrational: The Hidden Forces That Shape Our Decisions*. London, U.K.: Harper Collins, 2008.
8. D. Kahneman, P. Slovic, and A. Tversky, *Judgement Under Uncertainty: Heuristics and Biases*. Cambridge, U.K.: Cambridge Univ. Press, 1982.
9. D. Kahneman and A. Tversky, “Subjective probability: A judgement of representativeness,” *Cogn. Psychol.*, vol. 3, no. 3, pp. 430–454, Jul. 1972.
10. D. Kahneman and A. Tversky, “On the psychology of prediction,” *Psychol. Rev.*, vol. 80, no. 4, pp. 237–251, Jul. 1973.
11. E. Brunswick, “Representative design and probabilistic theory in a functional psychology,” *Psychol. Rev.*, vol. 62, no. 3, pp. 193–217, May 1955.
12. N. Karelaia and R. M. Hogarth, “Determinants of linear judgement: A meta-analysis of lens studies,” *Psychol. Bull.*, vol. 134, no. 3, pp. 404–426, May 2008.
13. D. G. Ullman and M. E. Doherty, “Two determinants of the diagnosis of hyperactivity: The child and the clinician,” *Adv. Dev. Behav. Pediatr.*, vol. 5, pp. 201–211, 1984.
14. P. Hoffman, P. Slovic, and L. Rorer, “An analysis of variance model for an assignment of configural cue utilization in clinical judgment,” *Psychol. Bull.*, vol. 69, no. 5, pp. 338–349, May 1968.
15. C. P. Bradley, “Can we avoid bias?” *Brit. Med. J.*, vol. 330, no. 7494, p. 784, Apr. 2005.
16. D. V. Budescu, S. Broomell, and H.-H. Por, “Improving communication of uncertainty in the reports of the intergovernmental panel on climate change,” *Psychol. Sci.*, vol. 20, no. 3, pp. 299–308, Mar. 2009.

895 17. S. S. Stevens, "On the theory of scales of measurement," *Science*,
896 vol. 103, no. 2684, pp. 677–680, Jun. 1946.
897 18. L. A. Cox, Jr., "What's wrong with risk matrices?" *Risk Anal.*,
898 vol. 28, no. 2, pp. 497–512, Apr. 2008.
899 19. M. Parker, R. Benson, and H. Trainor, *Information Economics:
900 Linking Business Performance to Information Technology*.
901 Englewood Cliffs, NJ: Prentice-Hall, 1988.
902 20. J. V. Neumann and O. Morgenstern, *Theory of Games and
903 Economic Behavior*. Princeton, NJ: Princeton Univ. Press, 1944.
904 21. J. S. Dyer, "Remarks on the analytic hierarchy process," *Manage.
905 Sci.*, vol. 36, no. 3, pp. 249–258, Mar. 1990.
906 22. J. S. Dyer, "A clarification of 'Remarks on the analytic hierarchy
907 process'," *Manage. Sci.*, vol. 36, no. 3, pp. 274–275, Mar. 1990.
908 23. R. D. Holder, "Some comment on the analytic hierarchy process,"
909 *J. Oper. Res. Soc.*, vol. 41, no. 11, pp. 1073–1076, Nov. 1990.
910 24. R. D. Holder, "Response to holder's comments on the analytic
911 hierarchy process: Response to the response," *J. Oper. Res. Soc.*,
912 vol. 42, no. 10, pp. 914–918, Oct. 1991.
913 25. S. Schenkman, "Inducement of nonexistent order by the analytic
914 hierarchy process," *Decis. Sci.*, vol. 28, no. 2, pp. 475–482,
915 Apr. 1997.
916 26. J. Pérez, "Some comments on Saaty's AHP," *Manage. Sci.*,
917 vol. 41, no. 6, pp. 1091–1095, Jun. 1995.
918 27. J. Pérez, J. Jimeno, and E. Mokotoff, "Another potential
919 shortcoming of AHP," *TOP*, vol. 14, no. 1, pp. 99–111, Jun. 2006.
920 28. P. J. H. Schoemaker and C. C. Waid, "An experimental
921 comparison of different approaches to determining weights in
922 additive utility models," *Manage. Sci.*, vol. 28, no. 2, pp. 182–196,
923 Feb. 1982.
924 29. C. Tsai, J. Klayman, and R. Hastie, "Effects of amount of
925 information on judgment accuracy and confidence," *Org. Behavior
926 Human Decis. Process.*, vol. 107, no. 2, pp. 97–105, Nov. 2008.
927 30. C. Heath and R. Gonzalez, "Interaction with others increases
928 decision confidence but not decision quality: Evidence against
929 information collection views of interactive decision making,"
930 *Org. Behavior Human Decis. Process.*, vol. 61, no. 3,
931 pp. 305–326, Mar. 1995.
932 31. M. L. Williams, A. R. Dennis, A. Stam, and J. E. Aronson,
933 "The impact of DSS use and information load on errors and
934 decision quality," *Eur. J. Oper. Res.*, vol. 176, no. 1, pp. 468–481,
935 Jan. 2007.
936 32. B. Flyvbjerg, *Procedures for Dealing With Optimism Bias in
937 Transport Planning*. London, U.K.: Brit. Dept. Transp., 2004.
938 33. A. H. Murphy and R. L. Winkler, "Can weather forecasters
939 formulate reliable forecasts of precipitation and temperature?"
940 *Nat. Weather Dig.*, vol. 2, no. 2, pp. 2–9, 1977.
941 34. D. Hubbard, *How to Measure Anything: Finding the Value of
942 Intangibles in Business*. Hoboken, NJ: Wiley, 2007.
943 35. S. Lichtenstein, B. Fischhoff, and L. D. Phillips, "Calibration of
944 probabilities: The state of the art to 1980," in *Judgement Under
945 Uncertainty: Heuristics and Biases*, D. Kahneman, P. Slovic, and
946 A. Tversky, Eds. Cambridge, U.K.: Cambridge Univ. Press,
947 1982, pp. 306–334.
948 36. J. M. Hampton, P. G. Moore, and H. Thomas, "Subjective
949 probability and its measurement," *J. R. Stat. Soc., Ser. A
950 (General)*, vol. 136, no. 1, pp. 21–42, 1973.
951 37. E. Servan-Schreiber, J. Wolfers, D. M. Pennock, and B. Galebach,
952 "Prediction markets: Does money matter?" *Electron. Markets*,
953 vol. 14, no. 3, pp. 243–251, Sep. 2004.
954 38. D. Bearden, C. Freaner, R. Bitten, and D. Emmons, "An
955 assessment of the inherent optimism in early conceptual designs
956 and its effect on cost and schedule growth," in *Proc. SSCAG/
957 SCAF/EACE Joint Int. Conf.*, Noordwijk, The Netherlands, 2008.
958 39. G. S. Simpson, F. E. Lamb, J. H. Finch, and N. C. Dinnie, "The
959 application of probabilistic and qualitative methods to asset
960 management decision making," presented at the Proc. SPA Asia
961 Pacific Conf. Integr. Modeling Asset Manage., Yokohama, Japan,
962 2000, Paper 59455.
963 40. W. Bailey, B. Couet, F. Lamb, G. Simpson, and P. Rose, "Taking a
964 calculated risk," *Oilfield Rev.*, vol. 12, no. 3, pp. 20–35, 2000.

Received March 9, 2009; accepted for publication
May 15, 2009

Douglas Hubbard *Hubbard Decision Research, Glen Ellyn,
IL 60137 USA (dwhubbard@hubbardresearch.com)*. Mr. Hubbard has
an M.B.A. degree in management information systems from the
University of South Dakota. He is the author of *How to Measure
Anything: Finding the Value of Intangibles in Business* (Wiley, 2007)
and *The Failure of Risk Management: Why It's Broken and How to
Fix It* (Wiley, 2009). He has published articles in *Information Week*,
CIO Enterprise, *Architecture Boston*, and *DBMS Magazine*. He was
formerly with Coopers & Lybrand and has more than 20 years of
experience in IT management consulting, including 14 years of
experience specifically in teaching organizations to use his AIE
(Applied Information Economics) method. His other professional
experience includes managing large IT strategy and software projects
in insurance, manufacturing, nuclear power, banking, and
pharmaceuticals.

Dylan Evans *University College Cork, School of Medicine,
Brookfield Health Sciences Campus, College Cork, Cork, Ireland
(d.evans@ucc.ie)*. Dr. Evans is Lecturer in Behavioral Science in the
School of Medicine at University College Cork, Ireland. After receiving
his Ph.D. degree in philosophy from the London School of Economics,
he conducted postdoctoral research in philosophy at King's College
London and in robotics at the University of Bath before moving to
the University of the West of England where he was Senior Lecturer in
Intelligent Autonomous Systems. In January 2008, he moved to the
University College Cork to become a Senior Research Scientist at the
Department of Computer Science, where he conducted research on
decision theory and risk management, before taking his current post.

AUTHOR QUERY

No query.