



MODEL RISK  
MANAGERS'  
INTERNATIONAL  
ASSOCIATION

# Request for Information and Comment: Financial Institutions' Use of Artificial Intelligence, including Machine Learning

---

May 21, 2021

Dr. Peter Quell

Dr. Joseph L. Breeden

# Summary

---

An increasing number of business decisions within the financial industry are made in whole or in part by machine learning applications. Since the application of these approaches in business decisions implies various forms of model risks, the Board of Governors of the Federal Reserve System, the Bureau of Consumer Financial Protection, the Federal Deposit Insurance Corporation, the National the Credit Union Administration, and Office of the Comptroller of the Currency issued an request for information and comment on the use of AI in the financial industry.

The Model Risk Managers' International Association MRMIA welcomes the opportunity to comment on the topics stated in the agencies' document. Our contact is [aimrm@mrmia.org](mailto:aimrm@mrmia.org).

# TABLE OF CONTENTS

<b>Summary .....</b>	<b>2</b>
<b>1. Introduction to MRMIA.....</b>	<b>4</b>
<b>2. Explainability.....</b>	<b>4</b>
<b>3. Risks from Broader or More Intensive Data Processing and Usage.....</b>	<b>7</b>
<b>4. Overfitting.....</b>	<b>9</b>
<b>5. Cybersecurity Risk .....</b>	<b>10</b>
<b>6. Dynamic Updating .....</b>	<b>11</b>
<b>7. AI Use by Community Institutions.....</b>	<b>12</b>
<b>8. Oversight of Third Parties .....</b>	<b>13</b>
<b>9. Fair Lending .....</b>	<b>13</b>
<b>10. Additional Considerations.....</b>	<b>16</b>
<b>11. Authors.....</b>	<b>20</b>

# 1. Introduction to MRMIA

---

Model Risk Managers' International Association (MRMIA.org) is a non-profit industry advocacy group focused on promoting best practices in model risk management. With a membership primarily from financial institutions, MRMIA seeks to provide an independent voice on model risk management that explicitly takes into account a cost-benefit analysis behind the practices being promoted. Specifically, MRMIA seeks to avoid chasing methods with diminishing returns, and instead seeks to promote those practices with real benefits to model developers, model owners, and those in various risk management functions who need to oversee the process.

In what follows, the questions from the RFI are addressed in sequential order.

## 2. Explainability

---

**Question 1:** How do financial institutions identify and manage risks relating to AI explainability? What barriers or challenges for explainability exist for developing, adopting, and managing AI?

Perhaps the most difficult issue when validating machine learning models is the explainability and interpretability of the model. This topic is being addressed at large in the industry, and some tentative industry standards are beginning to emerge:

- LIME (Local Interpretable Model-agnostic Explanations) analysis provides local interpretation of the machine learning algorithm's behavior using a linear model within some local region of feature space.
- SHAP (SHapley Additive exPlanations) is a way to measure and visualize the marginal contribution of each feature to the machine learning solution.
- Surrogate models use simple models, say, rule-based or decision trees, to explain the broad behavior of a machine learning algorithm.

Each of these have weaknesses: LIME and SHAP give only a partial view of the behavior of the machine learning algorithm across a local region or for one feature at a time, whereas surrogate modeling provides only an approximate, incomplete explanation.

Approximation methods such as LIME and SHAP may not be fully appropriate, especially for critical applications in finance where precise explanations are required. Some machine learning algorithms are intrinsically explainable. For example, decision tree builders such as ID3 or CART are a class of such algorithms. The decision tree that is constructed is both a non-linear model and is also immediately interpretable by a human analyst. More sophisticated models such as Deep ReLU networks can be interpreted exactly using equivalent sets of local linear models (LLM). With the appropriate toolkit, the analyst can determine how the network is behaving locally with as much precision as is required.

When considering explainability, it is worth considering carefully what is actually required:

- A full understanding of how the predictive algorithm works, or
- An explanation of how the algorithm makes a prediction or decision for one specific case at a time

The first problem is harder than the second and may be needed for model validation. However, the business may only require providing the second, e.g. being able to explain a lending decision to a specific customer.

**Question 2:** How do financial institutions use post-hoc methods to assist in evaluating conceptual soundness? How common are these methods? Are there limitations of these methods (whether to explain an AI approach's overall operation or to explain a specific prediction or categorization)? If so, please provide details on such limitations.

In terms of criteria for challenging the selection of a particular algorithm, model performance (predictive power) is certainly an important driver. The academic literature and some industry papers broadly coincide in the conclusion that machine learning algorithms (e.g. random forests, boosting, neural networks) generally outperform traditional models (e.g. logistic regressions, trees) when there are strong nonlinear relationships between the regressors and the target variable, whereas they show little to no improvement when the relationships are roughly linear. However, model performance should not

be the only criterion, and perhaps not even the most important one, when assessing the choice of an algorithm.

**Question 3:** For which uses of AI is lack of explainability more of a challenge? Please describe those challenges in detail. How do financial institutions account for and manage the varied challenges and risks posed by different uses?

AI is generally defined as data intensive because of its search for pockets of behavioral predictability. In measuring data, we need to look at breadth in number of factors and accounts and length in time. Data sets being used for machine learning are generally large in breadth but short in length. That means that the AI will not have context to understand that the patterns it learns apply to only the current environment. In the case of an environmental shift, such as just experienced with the COVID-19 pandemic, the AI approach will suffer.

Anti-fraud methods rely upon such short-in-time training sets, because the patterns are short in duration and change rapidly. Conversely, with lending and many other consumer relationships, the patterns play out over many years, so we cannot be so cavalier about focusing only upon recent data.

One technique proposed to deal with this situation uses an Age-Period-Cohort algorithm, that is optimized for long-range forecasting on long-in-time but thin-in-breadth data sets, to predict the mean of the distribution and then construct the machine learning method of choice to learn the behavioral patterns centered around that mean. Although this cannot account for all things lost by not having long-wide data sets, it does adjust for some of the worst aspects.<sup>1</sup>

---

<sup>1</sup>Breeden, J.L. and E. Leonova, "When Big Data Isn't Enough: Solving the long-range forecasting problem in supervised learning", *International Conference on Modeling, Simulation, Optimization and Numerical Techniques*, Shenzhen, China, 2019, DOI: 10.13140/RG.2.2.33311.05280.

# 3. Risks from Broader or More Intensive Data Processing and Usage

---

**Question 4:** How do financial institutions using AI manage risks related to data quality and data processing? How, if at all, have control processes or automated data quality routines changed to address the data quality needs of AI? How does risk management for alternative data compare to that of traditional data? Are there any barriers or challenges that data quality and data processing pose for developing, adopting, and managing AI? If so, please provide details on those barriers or challenges.

When validating machine learning models an initial hurdle is how to assess the data used for their development, including the transformations that are applied to them, either manually by model developers or automatically by machine learning techniques.

A first challenge is data representativeness. If data are gathered from a variety of sources, merged, transformed, and sometimes sampled and filtered using automated machine learning techniques, how can an independent validation verify that the resulting dataset is still representative of the population it intends to model? In this case, traditional representativeness analyses are usually still valid; regardless of the transformation process, an independent validation focuses on the resulting dataset and performs classical hypothesis tests (Ztest, Kolmogorov-Smirnov, chi-square, Kruskal-Wallis, etc.) to check for representativeness vis-à-vis the population to which the model will be applied, and any deviations from representativeness ought to be justified by model developers. It may be overly complex (and probably not worth it from a cost-benefit perspective) to assess every single step in the data transformation process when machine learning techniques are involved, and a result-oriented approach may be sufficient in the current state of the art.

A second challenge involves data traceability and data quality. How can an independent validation verify that each particular data point used in the model estimation process is correct and traces back to the source system with no errors in the course of the extract, transform and load (ETL) process, when data may be unstructured, sources may be diverse and machine learning techniques may have been applied in this process?

As for the assurance of the absence of errors in the ETL process, perhaps with machine learning techniques involved, a common approach is an audit-like data lineage control framework based on checkpoints. That is, a number of control checkpoints are established throughout the data transformation process, and simple comparison indicators are observed in each intermediate dataset (e.g. number of records, mean and variance of each variable at every checkpoint, etc.), together with the rationale provided by model developers about why and how records are transformed or filtered out. This is typically complemented with a sample-based traceability analysis of a number of individual records.

A third challenge involves feature engineering and the construction of synthetic variables that will be an input for the model estimation. How can an independent validation assess whether these synthetic variables are correctly built, pertinent and fit-for-purpose, especially when they have been created using machine learning techniques and may not necessarily have a business intuition of their own?

In this case, independent validation units are commonly focusing on two aspects. The first aspect is challenging the theoretical description of the algorithm that performs the feature engineering (when it is automated), whether it is conceptually sound, and if it has been correctly implemented (e.g. if a code library from a reputed source has been used). The second aspect is a deep analysis of the actual engineered variables in the resulting dataset. This includes both (1) qualitative questions, such as the business intuition of each variable or to what extent the engineered variable may be “contaminated” by the target variable (thus rendering it artificially predictive and flawed) and (2) quantitative checks such as the (linear or nonlinear) correlation with the target variable or the distribution (number and size) of the buckets in categorical variables.

**Question 5:** Are there specific uses of AI for which alternative data are particularly effective?

AI / ML techniques have been proven to have the flexibility needed for processing text and voice data in order to extract sentiment and general context. The applications of such analysis are still expanding.

Some alternate data sources are valuable, but with a particular need of AI / ML methods. Cash flow projection by analyzing deposit accounts, eBay accounts, etc. are examples of valuable data for consumer loans and small business loans, but the data is valuable without AI techniques.

Data sets involving device OS, time of day and day of week for loan applications might be analyzed with AI for fraud detection, but the same algorithms applied to loan underwriting carry too much risk of unintended bias.

## 4. Overfitting

---

**Question 6:** How do financial institutions manage AI risks relating to overfitting? What barriers or challenges, if any, does overfitting pose for developing, adopting, and managing AI? How do financial institutions develop their AI so that it will adapt to new and potentially different populations (outside of the test and training data)?

Because of the complexity of machine learning algorithms, they are prone to overfit training data, i.e. they fit spurious random aspects of the training data, as well as structure that is true for the entire population. There are techniques to handle overfitting, e.g. regularization, and all machine learning developers should be making use of these techniques.

However, just applying regularization is insufficient, and the performance of the machine learning solution should be measured to check for overfitting. Very simply, performance on the training data can be measured against that of an independent test data set. If it is very much better, then there is very likely an overfitting problem, and the machine learning developer should improve the machine learning algorithm to reduce the overfit, e.g. increase the regularization term, reduce model complexity and number of parameters, or increase sample size if possible.

The learning curve should also be used to demonstrate the performance of the learner against training data sample size. Essentially, the learning curve will show performance plotted against the amount of training data available. The performance can be measured both on training and test data. The primary objective is to see if there is enough data to train the learner (i.e., does performance converge to some upper bound) or if more data required. Secondly, since overfitting is related to sample size, it can also indicate how a machine learning algorithm is overfitting and can project how much data is required before overfitting becomes no longer relevant or reaches an acceptable level.

Many machine learning methods use performance on an out-of-sample data set to determine when to stop training the model. This approach to preventing overfitting is generally effective, but it carries a caveat. As a rule, the more often a data point is tested the less it can be considered out-of-sample. This was recognized several decades ago. In the case of hypothesis testing, the significance of the result should be adjusted based upon the number of tests conducted, as in the Holm–Bonferroni method.

When repeatedly testing scoring metrics or goodness-of-fit measures on an out-of-sample data set rather than hypothesis testing as above, the author is not aware of an equivalent adjustment, but the same principles apply. Simply stated, a good result with fewer out-of-sample tests is better than a slightly better result after many more tests. This needs to be considered when creating machine learning models and when reviewing completed work. One solution is to have a third and final out-of-sample test data set. Of course, this is only effective if it is used very few times.

These principles apply to both scoring models tested across hold out samples and time series models tested on an out-of-time sample. Rerunning an out-of-time test repeatedly can result in "look-ahead bias" where the meta-parameter decisions are based upon the analyst's judgement of accuracy on data that was supposed to be out-of-sample. This problem is particularly acute when modeling a short time series relative to the cycle being studied.

## 5. Cybersecurity Risk

---

**Question 7:** Have financial institutions identified particular cybersecurity risks or experienced such incidents with respect to AI? If so, what practices are financial institutions using to manage cybersecurity risks related to AI? Please describe any barriers or challenges to the use of AI associated with cybersecurity risks. Are there specific information security or cybersecurity controls that can be applied to AI?

In areas that use Deep Learning such as image processing or self-driving cars, researchers have found ways to fool machine learning algorithms with adversarial attacks. For example, change just a few pixels of an image to make a face recognition system fail to recognize a face, or add black sticky tape to a road sign to make a self-driving car speed illegally. This is a consequence of the fact that machine

learning algorithms are specialized, associative algorithms, rather than general intelligences with meaningful understanding of the world around them. The consequences for finance are uncertain, but we could imagine, for example, a fraudster fooling a machine learning system into providing credit with a careful selection of inputs.

The ability to fool machine learning algorithms relies on the attacker's ability to repeatedly test unusual boundary cases for which the AI may have been poorly trained and poorly understood by the developer. In these cases the attacker might subvert anti-fraud and credit risk controls. Intuitively, this is like extrapolating a polynomial model outside the range on which it was trained so that the outcomes can be unpredictable. To prevent these kinds of attacks, developers should carefully consider installing hard, simplistic boundaries around their AI framework that ignore unanticipated combinations of exotic factors and fall back simply to limits on such factors as bureau scores and loan-to-value ratios.

## 6. Dynamic Updating

---

**Question 8:** How do financial institutions manage AI risks relating to dynamic updating? Describe any barriers or challenges that may impede the use of AI that involve dynamic updating. How do financial institutions gain an understanding of whether AI approaches producing different outputs over time based on the same inputs are operating as intended?

Robustness in the statistical sense means that a statistical model still gives reliable parameter estimates even with deviations from modeling assumptions. Since machine learning algorithms work with few distributional assumptions, they are typically robust in this sense.

However, there is also a sense of robustness against changes in the distribution of the population, or population drift (sometimes also called concept drift). For the same reason that complex learning structures are prone to overfitting random variation in data, we might expect they are also prone to overfitting special conditions related to populations at a certain time. Thus, they may be more vulnerable to changes in population over time, in contrast to simpler linear models. In financial applications such as credit risk, population drift is typical. As a result, caution is required when developing machine learning solutions. Sensitivity to population drift can be tested by simulation study (artificially modifying

the distribution) or by backtesting over a long period of data. Again, if there is a problem, solutions such as regularization or simplifying the learning structure may help.

If a machine learning algorithm is not robust to population change it should not be deployed for stress testing, since it may be especially unrepresentative for extreme values that are expressed in a stress test scenario. In machine learning, there has been recent interest in Domain Adaptation Methods to address this problem of building a machine learning solution in one domain but applying it in a different domain.

This contribution is being written during the depths of the COVID-19 recession. As soon as shelter-in-place orders were issued in countries around the world, we all knew the corresponding credit risk models would have a problem. All the algorithms discussed here are data-driven pattern recognition engines. When past patterns are not predictive of future behavior, the models will fail to predict the outcome.

The advantage of machine learning algorithms may be described as the ability to find unique behavioral patterns that are not apparent to linear methods. However, behavioral patterns can shift suddenly, and time is required for new patterns to emerge from the data. Conversely, linear methods focusing on the borrower's cash flow and past payment performance are much more robust when conditions change. For these reasons, machine learning models should be viewed as more fragile than traditional financial modeling techniques. Since new data must be acquired for the ML algorithms to learn the new patterns, they can be unavailable for use during these transitions, so simpler traditional methods are probably required as a fallback.

## 7. AI Use by Community Institutions

---

**Question 9:** Do community institutions face particular challenges in developing, adopting, and using AI? If so, please provide detail about such challenges. What practices are employed to address those impediments or challenges?

Community banks and credit unions are unlikely to find much applicability of AI techniques to in-house data. To find patterns not visible to traditional techniques, the institution requires significantly more data

than is usually available. Vendor models may use AI / ML techniques as part of their solutions if they can leverage a multi-institution data set or public data like the mortgage data from Fannie Mae and Freddie Mac or securitized pool performance data.

Also, applying AI / ML techniques requires experience both in those techniques and in the oddities of financial institution data in order to avoid the risks of overfitting, unintended bias, and behavioral shifts. Few institutions possess such expertise at this time. This will change, but financial institutions need to recognize that their risk of failure is much higher than a machine learning algorithm that misunderstands a helpline voice request, so much greater caution is warranted for community institutions looking at AI techniques.

## 8. Oversight of Third Parties

---

**Question 10:** Please describe any particular challenges or impediments financial institutions face in using AI developed or provided by third parties and a description of how financial institutions manage the associated risks. Please provide detail on any challenges or impediments. How do those challenges or impediments vary by financial institution size and complexity?

Although vendors may be the best source of AI solutions for community lenders, those vendors cannot hide behind intellectual property claims when providing mission-critical solutions. Any system that provides anti-fraud, anti-money laundering (AML), or underwriting solutions must thoroughly test, document, and disclose their systems or they should not be considered for use. Often smaller lenders feel pressured to accept the documentation provided by large vendors, but this is a case where regulators can help force greater disclosure from vendors.

## 9. Fair Lending

---

**Question 11:** What techniques are available to facilitate or evaluate the compliance of AI-based credit determination approaches with fair lending laws or mitigate risks of non-compliance? Please explain

these techniques and their objectives, limitations of those techniques, and how those techniques relate to fair lending legal requirements.

With traditional linear methods, compliance with fair lending laws can be secured by excluding protected class status from the input data and generally restricting the training data to cash flow and payment performance metrics. With machine learning methods that incorporate alternate data, no such assurances are possible. No one can build a model that will not correlate to a variable that has not been observed. This is, by definition, impossible.

Tests that rely on the correlation of last names to ethnicity are weak and will only become weaker. Further, they say little about a range of discrimination risks. The only possible solution is for the government to change the laws around collecting protected class status so that lenders may test and prove that no such bias exists. Without data on which to test bias, model developers are blind. This is the single greatest obstacle to using machine learning on alternate datasets.

**Question 12:** What are the risks that AI can be biased and/or result in discrimination on prohibited bases? Are there effective ways to reduce risk of discrimination, whether during development, validation, revision, and/or use? What are some of the barriers to or limitations of those methods?

In our experience we have never met a lender that intended to be biased in lending decisions. The institutional risk is too great to warrant any perceived gain. The risk is of unintended bias, where a machine learning technique finds some unexpected and hard to explain patterns that give the lender an edge in pricing. Since lenders are not allowed to capture protected class data on loans other than for mortgages, lenders generally cannot perform tests to prove that their models are unbiased or rebuild them if they failed such a test. How is a lender to know that use of a certain device OS combined with purchasing patterns highly correlates to a protected class? In fact, such things are generally discovered in the news, when borrowers discovered the bias that lenders could not see.

These are cases where the machine learning method has done exactly what was asked, find behavioral patterns that correlate to higher risk. As a society we have chosen to classify some of those patterns as discriminatory. This can only be corrected by having data with which to impose constraints. For example, when building credit scores, lenders will often impose a constraint that those above a certain

age can be assigned no higher credit risk than those at the optimal age. This is possible because age is generally available. However, if gender, ethnicity, or other information is not in the data, then no such constraint is possible.

As an industry we need to stop thinking that blindness equates to being unbiased. This is certainly not true for machine learning.

**Question 13:** To what extent do model risk management principles and practices aid or inhibit evaluations of AI-based credit determination approaches for compliance with fair lending laws?

As mentioned above, the legal prohibition against collecting protected class status is the single reason that models cannot be built that explicitly comply with fair lending laws. To get unbiased models, regulators must allow for data collection. The only way lenders can find solutions to this problem is to find clever ways around the data collection prohibition, which amounts to gaming the system in ways that should not be encouraged.

**Question 14:** As part of their compliance management systems, financial institutions may conduct fair lending risk assessments by using models designed to evaluate fair lending risks (“fair lending risk assessment models”). What challenges, if any, do financial institutions face when applying internal model risk management principles and practices to the development, validation, or use of fair lending risk assessment models based on AI?

Fair lending risk assessment models are weak and will only become weaker. They are insufficient for constraining AI techniques and should not be used as a crutch.

**Question 15:** The Equal Credit Opportunity Act (ECOA), which is implemented by Regulation B, requires creditors to notify an applicant of the principal reasons for taking adverse action for credit or to provide an applicant a disclosure of the right to request those reasons. What approaches can be used to identify the reasons for taking adverse action on a credit application, when AI is employed?

Does Regulation B provide sufficient clarity for the statement of reasons for adverse action when AI is used? If not, please describe in detail any opportunities for clarity.

Reg B compliance is synonymous with the explainability problem described in Question 1. While only a few methods may have some global explainability, local explainability to the single applicant can usually be achieved. However, lenders need to consider how it will look if the answer relates to alternate data sources.

- “Your application was rejected because you subscribe to the following magazines: ...”
- “Your application was rejected because you applied from an off-brand Android tablet device.”
- “Your application was rejected because of the amount you spend monthly on alcohol.”

The requirement to make such disclosures would carry so much reputational risk to the lender that the use of “alternate data” must be carefully restricted to only those items that consumers would accept, even if they are in compliance with government regulations.

## 10. Additional Considerations

---

**Question 16:** To the extent not already discussed, please identify any additional uses of AI by financial institutions and any risk management challenges or other factors that may impede adoption and use of AI.

In the current state-of-the-art, there are at least three ways in which model validation teams in major banks can profit from AI:

- Challenger models
- Automated model building
- Optimization of the model validation process

The first way has already been discussed to some extent and consists in the development of machine learning models to challenge the model being validated. The challenger model is not necessarily intended to replace the current model (as mentioned above, it may encounter supervisory difficulties),

but it nonetheless makes sense to try machine learning challenger models, since they may reveal hidden relationships among regressors and non-linear behaviors that may be otherwise overlooked.

The second way consists in the utilization of tools that automatically produce a massive number of models based on the same input dataset, selects the ones that seem like a best fit, and offers them to the user for review. The concept of “best fit” may be based on many factors besides model performance, including compliance with a set of user restrictions or a well-balanced model in terms of weights of presence of variables from different profiles.

Again, even if this selection of models will not necessarily yield a replacement for the current model, it certainly provides useful insights such as an estimation of the maximum “predictability” capacity of the dataset (i.e., what is the highest performance that can ever be expected with this dataset, regardless of whether the model makes business sense?) and unexpected combinations of regressors that a person would probably not try.

The third way is probably the most difficult and the least explored: the automation or semi automation of the model validation process using machine learning techniques. This includes ideas such as using machine learning algorithms for cross-validation, feature extraction, and backtesting, but also more daring concepts like the drafting of the model validation report using natural language processing.

In the current state of the art, some progress has been made in this direction, but these initiatives are still in their very early stages. They are mostly limited to the automatic production of key performance indicators (quite similar to model monitoring) -- but not necessarily using machine learning techniques -- or to the automated filling of the test results in the periodic model validation report, reducing the manual workload, but again with no real intervention of advanced machine learning techniques.

**Question 17:** To the extent not already discussed, please identify any benefits or risks to financial institutions' customers or prospective customers from the use of AI by those financial institutions. Please provide any suggestions on how to maximize benefits or address any identified risks.

Currently model development processes were designed around traditional linear models. This introduces risks that are subtle, but significant.

First, anyone using a development package like R, SAS, Python, and others can quickly access and test the development of dozens of types of machine learning models and automatically search meta-

parameter space (though often slowly) to optimize so called “out-of-sample” performance. Problems with data being out-of-sample in such processes were addressed earlier. At the end of that process, what has the analyst learned?

Fraud detection, underwriting, pricing, and other mission critical systems are more than just models. A modeler who can see into the decision-making of the machine learning algorithm can relate those concepts to management for more systemic corrections.

- Conversation, Option 1:

Manager: “Why is your AML system rejecting all transactions in Idaho?”

Model Developer: “I don’t know. That’s what our vendor’s AML AI system says.”

Manager: “That’s not acceptable.”

- Conversation, Option 2:

Manager: “Why is your AML system rejecting all transactions in Idaho?”

Model Developer: “It appears that all of the transactions happening there are actually being rerouted from a prohibited off-shore location.”

Manager: “Then let’s talk to IT about filtering those transactions out of flow so that we can restore service to our customers in Idaho.”

We assume that Option 1 would not end there, but instead would launch an investigation. Better would be a proactive conversation.

- Conversation, Option 3:

Model Developer: “In developing our new model, we discovered that a high percentage of the AML alerts were coming via Idaho.”

Manager: Let’s talk to IT about filtering those transactions out of flow so that we can restore service to our customers in Idaho. Then maybe you can rebuild your model focusing on other causes.”

This is a trivial but obvious example. Insights related to the underlying causes of attrition, charge-off, or fraud can only be seen if the model developer can see into the model. Explainability is more than just an issue for validators, regulators, and compliance. Explainability speaks directly to institutional learning, or the lack thereof.

Second, machine learning introduces the risk of “validation arbitrage” or more specifically “p-Value arbitrage”. Because machine learning methods do not provide the same kinds of diagnostics as linear methods, they are already being held to a different standard than linear models. The most glaring example is in p-Values. A linear model can be rejected because the statistical significance of a coefficient does not meet some criterion. However, a simple neural net could replicate the same model, pass an out-of-sample test, and be put into production with no complaints about insignificant coefficients.

This approach to validation provides an advantage to black boxes. Many have complained that AI methods are at a disadvantage because of their black box nature, but if they closely examine how validations are performed, they will find that the opposite can be true. Of course, no modeling technique should have an advantage or disadvantage just because of how it is being reviewed in validation. Instead, we need a rethink of what is appropriate in validation. If we are to stop thinking linearly in model development, then we must also stop thinking linearly in validation, audit, and regulation.

# 11. Authors

---

## Peter Quell

Dr. Quell is Head of the Portfolio Analytics Team for Market and Credit Risk in the Risk Controlling Unit of DZ BANK AG in Frankfurt. He is responsible for methodological aspects of Internal Risk Models, Economic Capital and Model Risk. Prior to joining DZ BANK AG, Peter was Manager at d-fine GmbH where he dealt with various aspects of Risk Management Systems in the Banking Industry. He holds a MSc. in Mathematical Finance from Oxford University and a PhD in Mathematics. Peter is member of the editorial board of the Journal of Risk Model Validation and a founding board member of the Model Risk Management International Association ([mrmia.org](http://mrmia.org)).

## Joseph L. Breeden

Dr. Breeden has been designing and deploying risk management systems for loan portfolios since 1996. He founded Prescient Models in 2011, which focuses on portfolio and loan-level forecasting solutions for pricing, account management, CCAR, and CECL. He co-founded Deep Future Analytics in 2013 as a CUSO to bring solutions to credit unions and community banks. He is member of the board of directors of Upgrade, a San Francisco-based FinTech, an Associate Editor for the Journal of Risk Model Validation and for the Journal of Credit Risk, and a founding board member of the Model Risk Management International Association ([mrmia.org](http://mrmia.org)). Dr. Breeden received separate BS degrees in mathematics and physics in 1987 from Indiana University. He earned a Ph.D. in physics in 1991 from the University of Illinois studying real-world applications of chaos theory and genetic algorithms.