

THE THREE "RS" OF PREDICTIVE ANALYTICS

BUSINESS
CONSULTANTS

DEEP
TECHNOLOGISTS



*THE THREE "Rs" OF PREDICTIVE
ANALYTICS*

As companies commit to big data and data-driven decision making, the demand for predictive analytics has never been greater. While each day seems to bring another story of analytics success, it's important to keep some simple rules of thumb in mind to ensure that analytics solutions are both effective and adopted by your business users. Whether you're using traditional regression techniques or the latest machine-learning algorithms, these rules are critical to creating an analytics-driven enterprise.

West Monroe Partners calls these rules of thumb the three "Rs" of predictive analytics: reliable, repeatable, and relatable.



THE THREE "Rs" OF PREDICTIVE ANALYTICS

THE FIRST R: RELIABLE

It seems obvious, but in order to be used effectively, a predictive model needs to be accurate. Statisticians and modelers typically refer to this as "goodness of fit," or how well the model "fits" the data being modeled. There are many, many measures of goodness of fit, and the best statistic depends on the type of model you're building. For example, if you're forecasting sales for a future time period, Mean Absolute Percent Error (MAPE) is the classic measure. A MAPE of zero means that the model produces a perfect forecast, so a smaller MAPE indicates a more accurate model. For linear regression, the most common measure is the coefficient of determination, or adjusted R-squared, where a higher R-squared indicates that the model "explains" the data more accurately. For classification models such as logistic regression or CHAID/decision-trees – models that predict the likelihood of an event such as customer purchase or attrition – there is some debate since there are many different tests that tell you how well the model predicts the event's occurrence. With so many statistics at their disposal, it's not uncommon for statisticians to get trapped in an endless pursuit of goodness of fit rather than knowing when the model is "good enough" and that further gains in fit will have little business impact.

In order to avoid this, don't think of accuracy in absolute terms; that is, in terms of how often the model is "right." The truth is that most analytics these days are comparative in nature where outcomes are measured by their likelihood or rate relative to the average or random outcome. Most often, this is represented by an index or a lift calculation; for example:

- Wireless "Customer X" is 50-percent MORE likely to attrite than average based on her demographics, usage, and payment history
- A retail customer is more likely to purchase a particular product because his shopping pattern is similar to that of certain customers who purchase that product
- Sales are predicted to be 25-percent higher in the Midwest than in the Southeast

In all of these cases, the actual value of the prediction or forecast is less important than the relative lift or improvement the model provides compared to having no model at all.

When used for comparative purposes, a high degree of model precision simply isn't necessary. Instead, your goal should be to build the simplest model possible that produces the most consistent results. A model that is quick to build and easy to deploy will start producing business gains immediately. Likewise, even models with relatively low lift or that perform only marginally better than random can produce very high return on investment at scale. Eric Seigel, author of *Predictive Analytics: The Power to Predict Who Will Click, Buy, Lie, or Die*, calls this the "prediction effect" where "a little prediction goes a long way." Think of an e-commerce retailer that is dealing with thousands of customers each day through its website. Improving cross-sell by only a percentage point or two through model-driven recommendations during check out can produce very high returns in a matter of weeks when applied across all check outs. This is also true for applications such



THE THREE "Rs" OF PREDICTIVE ANALYTICS

as customer churn. Imagine a cable company whose average customer generates \$1,000 to \$1,500 in annual revenue and that loses 20 percent of these customers to its competition annually. Even small improvements in customer retention can lead to big gains for the company, especially since these gains compound over time as customer retention increases.

In both cases, simple, effective models can pay for themselves many times over, even if they aren't theoretically the "best" models possible.

Another danger of pursuing model accuracy, or goodness of fit, is called "over-fitting." Simply put, this refers to cases where a model fits the data that it's built on so well that it won't hold up with new data. The model "over-fits" the data at hand and is simply too good to be true. As modeling becomes more automated through machine-learning and applied to "big data," the risk of over-fitting increases greatly because the algorithms themselves are trained to maximize goodness of fit, and even very small correlations in the data become statistically significant. Models that are over-fit are overly complex and are more susceptible to small changes in the underlying data.

In order to avoid this, an analyst should always test models on hold-out samples of different customers or time periods and compare model performance. Disciplined cross-validation of models will often lead you to a model that may have lower goodness of fit during analysis but that produces better results over time.

Finally, as described in our first big data white paper, [*Could vs. Should: Balancing 'Big Data' and Analytics Technology with Practical Outcomes*](#), you should always balance your quest for model accuracy against the cost of being wrong. In the e-commerce case above, the incremental cost of making a wrong product recommendation is very low. If your model is forecasting next quarter's product demand in order to set manufacturing and inventory levels, however, you'll want to understand what the business costs are in terms of over or under supply, so again you know when your model is "good enough" rather than the best possible model.

Like many pursuits, the law of diminishing returns will eventually take hold. In the world of predictive modeling, there will be a point when a higher degree of statistical accuracy simply won't change the business outcome or return on investment. There is no such thing as a perfect model, so rather than trying to build one, use a mix of statistics and business understanding to know when you've built an actionable and effective predictive model.

THE SECOND "R": REPEATABLE

There are two facets to this "R" – repeatable results and repeatable process.

Related to the first "R" above, a model has to be able to reproduce results across customers, markets, time



THE THREE "Rs" OF PREDICTIVE ANALYTICS

periods, etc. A model that is over-fit can't be relied upon to produce good results repeatedly. Again, cross-validation during model development helps ensure that a model's results can be repeated over time. You're not done once the model is deployed, however. It's also critical that the model be continually tested and measured on small, random control samples to ensure that it still provides the lift in performance expected when it was built. These control samples will help you identify when a model is in need of repair or rebuild, while also providing built-in samples for new model development.

Building a repeatable modeling process will help you create a truly analytics-driven enterprise. Rather than viewing modeling as one-off, episodic projects, you should be building a process that can be repeated on new products or business problems. This is often the last thing on the modeler's mind as he or she is in the throes of model building; however, building a repeatable data collection process and modeling methodology will create true business advantage, while also increasing the overall efficiency of your modeling team. Clear data documentation, well-documented code, and standard business metrics and deliverables will ensure you produce a modeling "practice" rather than simply good models.

Fortunately, the Cross-Industry Standard Process for Data Mining (CRISP-DM) provides an excellent reference for building a repeatable modeling process. The diagram on the next page provides a view of the CRISP-DM process, which offers a great way to organize your methods and procedures to build a repeatable modeling process.



THE THREE "Rs" OF PREDICTIVE ANALYTICS

CROSS-INDUSTRY STANDARD PROCESS FOR DATA MINING (CRISP-DM)



Business Understanding

- Design and document business requirements, constraints, and objectives
- Produce a project plan with the timelines and tools to be used

Data Understanding

- Collect, integrate, and study data

Data Preparation

- Cleanse data, sample data, develop key performance indicators, and derive attributes for analysis

Modeling

- Select and apply modeling techniques to build predictive or descriptive model
- Test and validate the model

Evaluation

- Produce final report illustrating performance and key business insights
- Review or revise based on stakeholder feedback
- Conduct final project review

Deployment

- Implement ongoing scoring, planning, monitoring, and maintenance

Source: [IBM](#)

THE THIRD "R": RELATABLE

The final "R" refers to the ability of your business clients to "relate" to the predictive model. This is often the most important, yet least appreciated, of the three "Rs" when it comes to actually putting a predictive model into practice. Business users must be able to understand the predictive model's insights and performance in terms they can relate to. Only when the analytics are presented in a way that is meaningful to the business will the user trust the model enough to build a business strategy around it.

Unfortunately, modelers often put this "R" last in their priorities. Instead, they focus on the first "R," reliability, and demonstrate the model's accuracy in statistical rather than business terms. Or they focus on building an elegant or theoretically interesting model, rather than a simple, effective model that produces business results. Modelers who don't recognize this are often frustrated when their models aren't deployed. The bottom line is that no matter how good a model is, if the user can't relate to it, the model will be either



THE THREE "Rs" OF PREDICTIVE ANALYTICS

ignored or underutilized.

Two short examples illustrate this point.

In the first example, the head of analytics for a Fortune® 500 entertainment company was frustrated by his chief marketing officer's desire to "understand" his models. The analytics leader lamented that his boss will "never understand" his model, since he wasn't a trained expert in modeling and statistics. He told the CMO to "trust him" that the model will work. What he failed to realize was that the CMO wasn't seeking to understand the model from a technical point of view. Instead, the CMO wanted to know enough about the model so he could trust it in his business—he wanted to believe the model, not understand it. The analytics leader's failure to make his model relatable and to present it in business rather than statistical terms led to the model sitting on the shelf to this day, limiting his organization's adoption of predictive modeling in general.

In contrast, the second example provides a successful illustration of making a model relatable. In this example, a globally-integrated supply chain company was being challenged by its client on the accuracy of its forecasts. The client is a leading global fast-food retailer and was relying on its supply chain partner for daily and weekly demand forecasts to optimize inventory at each restaurant. The client used these forecasts to generate recommended orders for each restaurant every morning, and it suspected that the model's forecast accuracy was decreasing. In the past, the modeler would typically address this concern by comparing MAPE over time to show that the forecasting model is still performing well and producing a low error rate on average. In this case, however, the modeler took a different, much more effective approach. The modeler took different items, such as hamburger patties, potatoes, etc., and translated the average error rates into actual orders or units of inventory. In other words, if patties can only be purchased in packages of 10 with 10 packages per case, he calculated the percentage of time when the forecast was off enough to order more or fewer cases. By defining the model's performance in the language of the restaurant manager, rather than the language of the statistician, the modeler was able to show that even if the actual forecasts of demand have become less accurate, it had virtually no effect on the actual orders placed at the restaurant, since the error was rarely large enough to change the number of cases ordered. As a result, the models remain in use and the restaurant managers, themselves, have now become advocates for the modeling solution in the field.

Predictive analytics deliverables or presentations often illustrate the divide between the desire of the modeler to "show his work" and the need of the business user to trust the model. When left to his or her own devices, the modeler will often devote the bulk of content to the data that was collected and studied, how it was prepared, and the modeling methodology used. Model results, insights, and recommended uses are often saved for the closing sections of the report. While this approach may satisfy the modeler, it is the opposite of what the business user needs. In order to encourage the adoption of predictive analytics, the modeler should reverse his or her approach. The most effective modeling presentations BEGIN with the



THE THREE "Rs" OF PREDICTIVE ANALYTICS

modeling challenge, then present the resulting model's performance and key insights in business terms. The report should end with recommendations for model deployment and use. Data documentation, modeling methodology, and summary statistics should be in an appendix or an accompanying document for those who require it. This doesn't mean that those elements are not important; it simply means that they shouldn't be the focus of the model discussion.

The challenge with this third "R" is that it is rarely taught as part of any analytics curriculum. Like the analytics leader above, modelers typically learn this the hard way as they watch perfectly good models go unused because the business "didn't get it."

While the world of big data and predictive analytics seems to be ever-changing, it's important to keep these simple "three Rs" in mind when building an analytic enterprise. By focusing on making predictive analytics reliable, repeatable, AND relatable, you will be able to turn your business insight into foresight that drives business growth.