



Combining Models & Ensembles

2016 CAS Ratemaking and Product Management Seminar

Christopher Cooksey, FCAS, MAAA

Chief Actuary, EagleEye Analytics

AGENDA

1

Styles of combining models

2

Ensembles

3

Objections to ensembles

What do people
mean when
they say
they are building
"a GLM"?

- ▶ Solving for the optimal
relativities?
- ▶ Determining the best
structure among the possible
predictors?
- ▶ Or more generally using a
GLM technique to solve a
business problem?

Consider using a GLM technique to build a personal lines auto rating algorithm

▶ What do you
do with VIN?

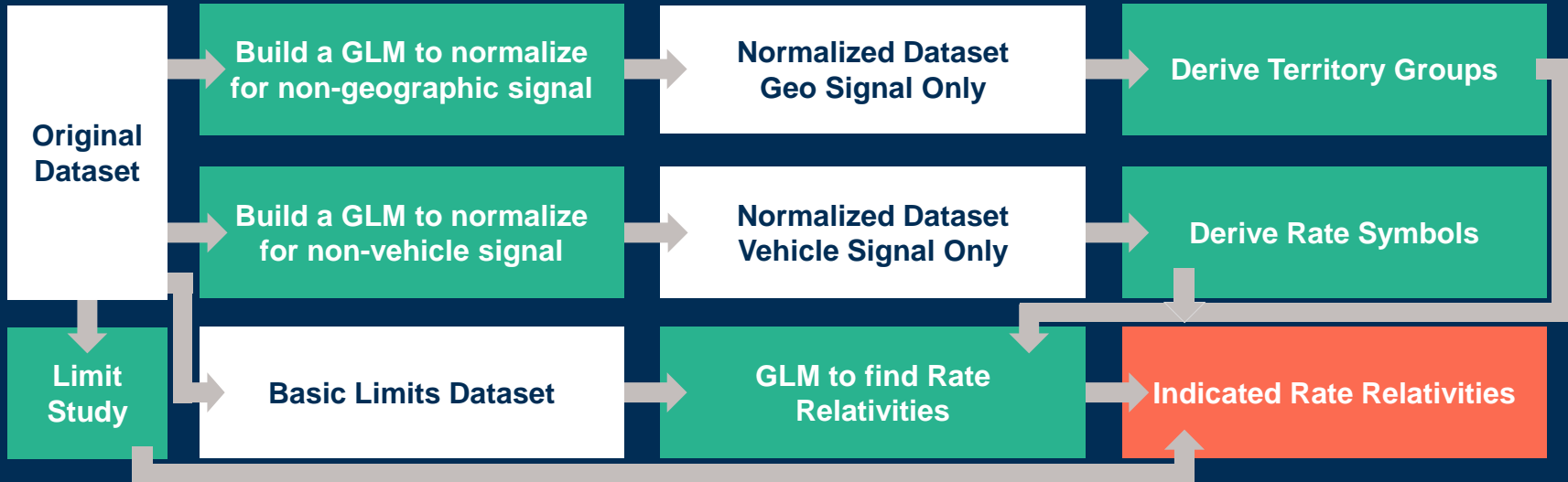
▶ How do you
capture
geographic
signal?

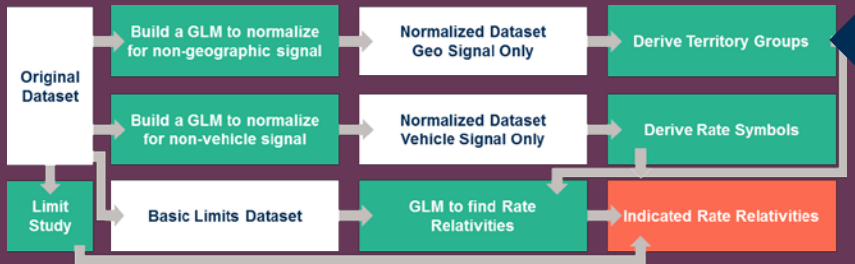
▶ Do you model
limit directly?

How to approach the geographic portion of the signal?

- ▶ A separate study on the geographic signal is a fine idea, but doesn't the data include correlations with non-geographic predictors?
- ▶ Often we build a GLM for the sole purpose of normalizing the data for everything BUT the geographic signal. This gives us the necessary data to do a geographic study.

Building “a GLM” for an auto rating plan may look like this...





The previous auto example contains a nested model

Sometimes nested models are built to represent specific portions of the signal, but sometimes they are motivated by the need to stabilize the model.

- ▶ Loss ratio can be a difficult target; it is inherently a residual metric.
- ▶ Building models which target frequency and severity can sometimes help.
- ▶ Like territory, the output of the sub-models are used as predictors in the desired model.

Application of predictive analytics in insurance is much wider than rating

Claims
Example:
**Is attorney
involvement
predictive of
claim severity?**



**For claims where it is
too early for an
attorney, would a model
predicting attorney
involvement improve
the model predicting
severity?**



Here the model
is trying to fill in
the blank for
what we know
is an important
piece of
information.

The signal in the data can be compartmentalized
(A key concept we've been using but not discussing)

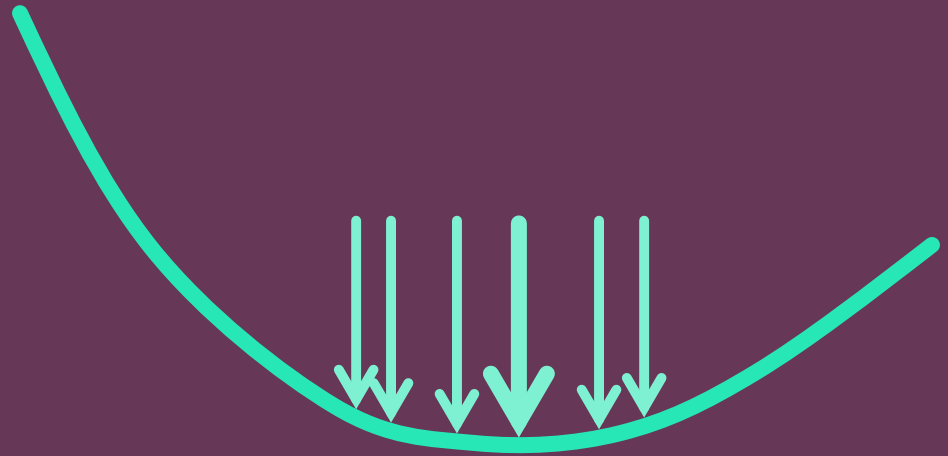
The “geographic portion” of the signal, for example:

- ▶ Keeping track of what parts of the signal (information) in the data has been represented is a critical part of modeling.
- ▶ A GLM is by its nature a linear model – **how about the non-linear portion of the information in the data?**

Another Twist

If you have two models, each of which perform similarly from a statistical perspective, which do you choose?

Normally we work with some function to define “best.”



Multiplicity of Models

...there is often a multitude of different descriptions [equations $f(x)$] in a class of functions giving about the same minimum error rate.

Breiman, L. (2001). Statistical Modeling: The Two Cultures. *Statistical Science*, Vol. 16, No. 3.

Data will often point with almost equal emphasis on several possible models, and it is important that the statistician recognize and accept this.

McCullagh, P. and Nelder, J. (1989). *Generalized Linear Models*.

AN UNREALISTIC ILLUSTRATION

Ground Rules

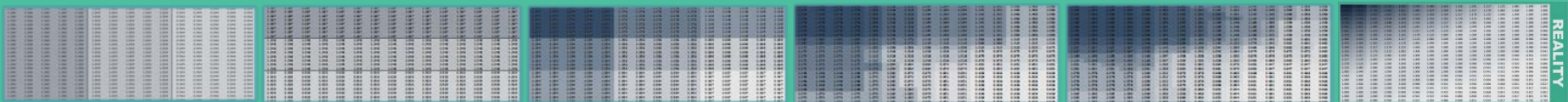
- 1) We get to know reality & compare our models directly.**
- 2) Assume the numbers are frequency relativities.**
- 3) Volume is limited; we can only divide the data into three equally-sized groups.**
- 4) Model predictions are just the average for each defined group.**

AN UNREALISTIC ILLUSTRATION

2.026	1.948	1.801	1.732	1.665	1.539	1.480	1.423	1.316	1.265	1.217	1.125	1.082	1.040	1.000
1.948	1.873	1.732	1.665	1.601	1.480	1.423	1.369	1.265	1.217	1.170	1.082	1.040	1.000	1.000
1.873	1.801	1.665	1.601	1.539	1.423	1.369	1.316	1.217	1.170	1.125	1.040	1.000	1.000	1.000
1.801	1.732	1.601	1.539	1.480	1.369	1.316	1.265	1.170	1.125	1.082	1.000	1.000	1.000	1.000
1.732	1.665	1.539	1.480	1.423	1.316	1.265	1.217	1.125	1.082	1.040	1.000	1.000	1.000	1.000
1.665	1.601	1.480	1.423	1.369	1.265	1.217	1.170	1.082	1.040	1.000	1.000	1.000	1.000	1.000
1.601	1.539	1.423	1.369	1.316	1.217	1.170	1.125	1.040	1.000	1.000	1.000	1.000	1.000	0.980
1.539	1.480	1.369	1.316	1.265	1.170	1.125	1.082	1.000	1.000	1.000	1.000	1.000	1.000	0.980
1.480	1.423	1.316	1.265	1.217	1.125	1.082	1.040	1.000	1.000	1.000	1.000	1.000	0.980	0.960
1.423	1.369	1.265	1.217	1.170	1.082	1.040	1.000	1.000	1.000	1.000	1.000	1.000	0.980	0.960
1.369	1.316	1.217	1.170	1.125	1.040	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.980	0.960
1.316	1.265	1.170	1.125	1.082	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.980	0.960	0.941
1.265	1.217	1.125	1.082	1.040	1.000	1.000	1.000	1.000	1.000	1.000	0.980	0.960	0.941	0.922
1.217	1.170	1.082	1.040	1.000	1.000	1.000	1.000	1.000	1.000	0.980	0.960	0.941	0.922	0.904
1.170	1.125	1.040	1.000	1.000	1.000	1.000	1.000	1.000	0.980	0.960	0.941	0.922	0.904	0.886
1.125	1.082	1.000	1.000	1.000	1.000	1.000	1.000	0.980	0.960	0.941	0.922	0.904	0.886	0.868
1.082	1.040	1.000	1.000	1.000	1.000	1.000	1.000	0.980	0.960	0.941	0.922	0.904	0.886	0.868
1.040	1.000	1.000	1.000	1.000	1.000	1.000	0.980	0.960	0.941	0.922	0.904	0.886	0.868	0.851
1.000	1.000	1.000	1.000	1.000	1.000	0.980	0.960	0.941	0.922	0.904	0.886	0.868	0.851	0.834
1.000	1.000	1.000	1.000	1.000	0.980	0.960	0.941	0.922	0.904	0.886	0.868	0.851	0.834	0.817
1.000	1.000	1.000	1.000	0.980	0.960	0.941	0.922	0.904	0.886	0.868	0.851	0.834	0.817	0.801
1.000	1.000	1.000	1.000	0.980	0.960	0.941	0.922	0.904	0.886	0.868	0.851	0.834	0.817	0.801
1.000	1.000	1.000	0.980	0.960	0.941	0.922	0.904	0.886	0.868	0.851	0.834	0.817	0.801	0.785
1.000	1.000	0.980	0.960	0.941	0.922	0.904	0.886	0.868	0.851	0.834	0.817	0.801	0.785	0.769

REALITY

AN UNREALISTIC ILLUSTRATION



Ensembles remain robust even as they become increasingly complex.

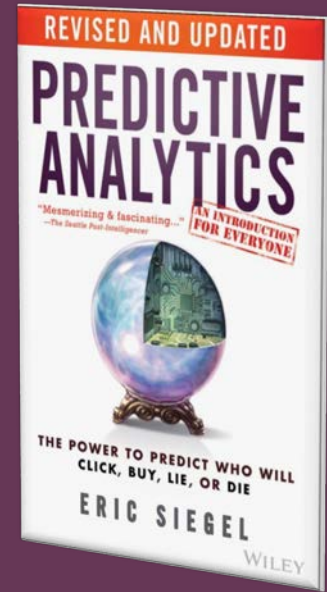
They seem to be immune to this limitation, as if soaked in a magic potion against overlearning.

Siegel, E. (2013). *Predictive Analytics*.

Ensemble modeling has taken the [Predictive Analytics] industry by storm.

It's often considered the most important predictive modeling advancement of this century's first decade.

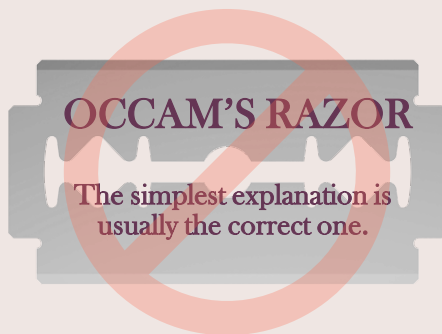
Siegel, E. (2013). *Predictive Analytics*.



OBJECTIONS TO ENSEMBLES

Resistance usually centers around complexity

Simpler is preferred in the absence of certainty, when multiple models perform equally well.



But if an ensemble performs better, then it is simply the better model.

OBJECTIONS TO ENSEMBLES

Framing the question as the choice between accuracy and interpretability is an incorrect interpretation of what the goal of a statistical analysis is.

The point of a model is to get useful information about the relation between the response and predictor variables. Interpretability is a way of getting information.

Breiman, L. (2001). Statistical Modeling: The Two Cultures. *Statistical Science*, Vol. 16, No. 3.

OBJECTIONS TO ENSEMBLES

(Consider neural
nets vs. trees)

All machine learning techniques are
equally difficult to explain.

Departments of insurance won't
accept them.

Because it can't be explained in simple terms,
there is no opportunity for insight.

Anything that is
theoretically
possible will be
achieved in
practice,
no matter what the
technical
difficulties are, if it
is desired greatly
enough.

~ Arthur C Clarke ~



OBJECTIONS TO ENSEMBLES

Don't think a complex model will be accepted for pricing in your underwriting-driven culture?

Context & Needs for Predictive Analytics in Insurance

- ✓ **Underwriting**
- ✓ **Marketing**
- ✓ **Claims management**
- ✓ **Internal monitoring**

- 1) Be aware of your language. Not everyone knows what you mean when you say you are going to “build a GLM”. Don’t short-sell the effort.
- 2) Combining models is typical. Expect it in any project are doing.
- 3) Keep track of the signal (information) in your data. Consider what you’ve done to represent relevant portions of it.
- 4) If your approach is linear, are there ways to capture the non-linear parts of your signal?
- 5) Results from ensemble approaches are transforming other industries and are worth the effort for insurance predictive modelers to explore.
- 6) The difficulties around model complexity are real, but can be addressed.

TAKEAWAYS



Chris Cooksey
Chief Actuary

ccooksey@EEAnalytics.com
855.757.8500

EEAnalytics.com

