Statistical Models and Credibility

Leigh J. Halliwell, FCAS, MAAA

Milliman & Robertson, Inc. 3 Garret Mountain Plaza West Patterson, NJ 07424

Casualty Actuarial Society 1998 Seminar on Ratemaking Chicago Hilton and Towers Friday, March 13, 1998

Outline

- 1. Matrices as Linear Mappings
- 2. The Linear Statistical Model
- 3. Credibility and Prior Information
- 4. Credibility and the Random-Effects Model
- 5. Conclusion

All this in seventy-five minutes?

1. Matrices as Linear Mappings

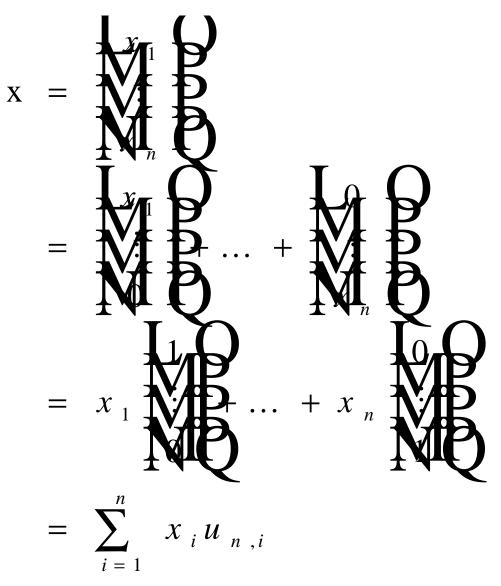
Vector
$$\mathbf{x} = \mathbf{x} \in \mathbb{R}^n$$

x is a point in *n*-dimensional real-number space. It packages *n* pieces of information.

How to multiply a vector by a scalar and how to add two *n*-dimensional vectors are obvious.

Define the unit vector $u_{n,i}$ as the R^n vector whose i^{th} element is one, the other elements being zeroes.

By the properties of addition and scalar multiplication,



Consider a linear mapping A: $R^n \rightarrow R^m$ Linear as to scalar multiplication: $A(\alpha x) = \alpha A(x)$ Linear as to vector addition: $A(x_1+x_2) = A(x_1)+A(x_2)$. In general,

Abg A
$$f_{i=1}^{x}$$
 A $f_{i=1}^{x}$ $x_{i}u_{n,i}$ $f_{i=1}^{n}$ $x_{i}AG_{n,i}h$

Therefore, a linear mapping is uniquely determined by to where it maps the unit vectors of R^n .

Let $A_i = A(u_{n,i}) \in \mathbb{R}^m$. Every linear mapping can be represented by the $m \times n$ matrix $[A_1 \dots A_n]$. The *i*th column of the matrix specifies the vector of \mathbb{R}^m to which A maps the *i*th unit vector of \mathbb{R}^n .

So,

$$A b g = \sum_{i=1}^{n} x_i A G_{n,i} h = \sum_{i=1}^{n} x_i A_i = \begin{bmatrix} A_1 & \cdots & A_n \end{bmatrix}$$

This looks like matrix multiplication, although matrix multiplication has not yet been defined (see slide 14).

How to multiply a matrix by a scalar and how to add two $m \times n$ matrices are obvious.

If A and B are two linear mappings from R^n to R^m , then,

$$\mathbf{b} + \mathbf{B} \mathbf{g} \mathbf{g} = [\mathbf{b} + \mathbf{B} \mathbf{g} \cdots \mathbf{b} + \mathbf{B} \mathbf{g}] \mathbf{b} \mathbf{g}$$
$$= \sum_{i=1}^{n} x_i \mathbf{b} + \mathbf{B} \mathbf{g}$$
$$= \sum_{i=1}^{n} x_i \mathbf{b} + \mathbf{B} \mathbf{g}$$
$$= \sum_{i=1}^{n} x_i \mathbf{b} + \mathbf{B} \mathbf{g}$$
$$= \sum_{i=1}^{n} x_i \mathbf{A}_i + \sum_{i=1}^{n} x_i \mathbf{B}_i$$
$$= \mathbf{A} \mathbf{b} \mathbf{g} + \mathbf{B} \mathbf{b} \mathbf{g}$$

Matrix addition is commutative and associative. There is a zero matrix, and every matrix has an additive inverse. These are the addition characteristics of rings.

But what about matrix multiplication?

Let B $(l \ m)$ represent a mapping from R^m to R^l . Let '•' represent the composition of mappings. B•A maps from R^n to R^l . But:

b • A
$$\mathfrak{G} \mathfrak{G} \mathfrak{G} = B \mathfrak{G} \mathfrak{A} \mathfrak{B} \mathfrak{G}$$

$$= B \mathfrak{G} \mathfrak{G} \mathfrak{A} \mathfrak{B} \mathfrak{G}$$

$$= B \mathfrak{G} \mathfrak{A} \mathfrak{B} \mathfrak{G} \mathfrak{G}$$

$$= \sum_{i=1}^{n} x_{i} B \mathfrak{B} \mathfrak{G}_{i} \mathfrak{G}$$

$$= [B \mathfrak{B} \mathfrak{A}_{1} \mathfrak{G} \cdots B \mathfrak{B} \mathfrak{A}_{n} \mathfrak{G} \mathfrak{B} \mathfrak{G}]$$

The columns of B•A make sense. For example, A maps $u_{n,1}$ to A_1 , then B maps A_1 to $B(A_1)$. So the *i*th column of B•A shows to what vector of R^l B•A maps $u_{n,i}$.

Composition (•) is always a commutative operator:

$$\mathbf{C} \cdot \mathbf{b} \cdot \mathbf{A} \underbrace{\mathbf{gh}}_{\mathbf{F}} = \mathbf{C} \underbrace{\mathbf{ch}}_{\mathbf{F}} \cdot \mathbf{A} \underbrace{\mathbf{gh}}_{\mathbf{F}} \underbrace{\mathbf{gh}}_{\mathbf{F}} = \mathbf{C} \underbrace{\mathbf{ch}}_{\mathbf{F}} \cdot \mathbf{A} \underbrace{\mathbf{gh}}_{\mathbf{F}} \underbrace{\mathbf{gh}}_{\mathbf{F}} = \mathbf{c} \underbrace{\mathbf{ch}}_{\mathbf{F}} \cdot \mathbf{B} \underbrace{\mathbf{gh}}_{\mathbf{F}} \underbrace{\mathbf{gh}}_{\mathbf{F}} \underbrace{\mathbf{gh}}_{\mathbf{F}} = \mathbf{b} \cdot \mathbf{B} \underbrace{\mathbf{gh}}_{\mathbf{F}} \underbrace{\mathbf{gh}} \underbrace{\mathbf{gh}}_{\mathbf{F}} \underbrace{\mathbf{gh}}_{\mathbf{F}} \underbrace{\mathbf{gh}}_{\mathbf{F}} \underbrace{\mathbf{gh}}$$

Composition distributes over matrix addition:

$$\mathbf{A} \cdot \mathbf{b} + \mathbf{C} \mathbf{g} \mathbf{b} \mathbf{g} = \mathbf{A} \mathbf{c} \mathbf{b} + \mathbf{C} \mathbf{g} \mathbf{g} \mathbf{g}$$

$$= \mathbf{A} \mathbf{G} \mathbf{b} \mathbf{g} + \mathbf{C} \mathbf{b} \mathbf{g} \mathbf{g}$$

$$= \mathbf{A} \mathbf{G} \mathbf{b} \mathbf{g} + \mathbf{A} \mathbf{c} \mathbf{c} \mathbf{b} \mathbf{g}$$

$$= \mathbf{b} \mathbf{a} \cdot \mathbf{b} \mathbf{g} \mathbf{g} + \mathbf{b} \cdot \mathbf{c} \mathbf{g} \mathbf{g}$$

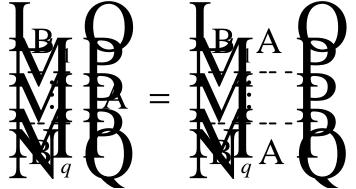
Being commutative and distributive, the composition of linear mappings behaves like multiplication. More accurately, multiplication of two matrices is really the composition of two mappings.

It is really easier to think of an $m \times n$ matrix A as the linear mapping A: $R^n \rightarrow R^m$. The columns of A show to where the unit vectors of R^n are mapped.

Matrices are linear mappings.

Partitioned mappings (matrices) compose (multiply) as follows:

b • A $\mathfrak{G} \mathfrak{G} = \begin{bmatrix} B \mathfrak{b}_{1} \mathfrak{g} \cdots B \mathfrak{b}_{n} \mathfrak{g} \mathfrak{b}_{n} \mathfrak{g} \end{bmatrix}$ B $\begin{bmatrix} A_{1} & \cdots & A_{p} \end{bmatrix} = \begin{bmatrix} B A_{1} & \cdots & B A_{p} \end{bmatrix}$



Recall from slide 6:

$$A b g = \sum_{i=1}^{n} x_i A \mathcal{G}_{n,i} h = \sum_{i=1}^{n} x_i A_i = \begin{bmatrix} A_1 & \cdots & A_n \end{bmatrix}$$

So,
$$\begin{bmatrix} A_1 & \cdots & A_p \end{bmatrix} = \sum_{i=1}^{p} A_i b_i g$$

In general,
$$\begin{bmatrix} B_1 & \cdots & B_p \end{bmatrix} = \begin{bmatrix} A_p \\ A_p \end{bmatrix} = \begin{bmatrix} B_p \\ A_i \end{bmatrix}$$

Combining the partition rules of this and the previous slide...

As long as the B and A are partitioned conformably, the *ij*th cell of the mapping B•A, or of the matrix product BA, will be :

$$\sum_{k} \mathbf{B}_{ik} \bullet \mathbf{A}_{kj}, \text{ or } \sum_{k} \mathbf{B}_{ik} \mathbf{A}_{kj}$$

Partitionwise multiplication is no different from elementwise multiplication. In fact, the elements *are* just the finest partitions, (1×1) partitions.

Matrix multiplication must have been first defined (by Cayley, Hamilton, Sylvester?) according to the interpretation of matrices as linear mappings. 'B times A' is 'B of A.'

2. The Linear Statistical Model

Information is an *n*-dimensional vector (*n*-tuple).

Models explain the more complicated by the less complicated.

Example: A car moves in a straight line at a constant velocity. At times t_1, \ldots, t_n its position is observed to be d_1, \ldots, d_n . Time is the independent variable and distance is the dependent variable. We define our time and distance scales such that at time zero the car is at a distance zero.

We know that there will be some number *r* such that $d_i = r \cdot t_i$.

The linear model for this example is:

$$d_{\mathbf{b}\times 1}g = \frac{d}{d\mathbf{b}} \frac{d}{d\mathbf{b}} \frac{d}{d\mathbf{b}} \frac{d}{d\mathbf{b}} \frac{d}{d\mathbf{b}} = t_{\mathbf{b}\times 1}g^{\mathbf{b}\times 1}g$$

If we know t, the independent variable, we can explain the ndimensional d in terms of the 1-dimensional r. Reducing complexity is the essence of modeling.

Predictive power is a by-product. When we know r, we can predict distances for new t_i s.

Here is a model which explains a *t*-dimensional phenomenon in *k*-dimensions (hopefully k < t):

$$y_{b \times 1}g = f \Theta_{b \times 1}g$$

f is some map from R^k to R^t . If *f* is a *linear* map, then we can express *f* as some $t \times k$ matrix X, and the model becomes:

$$y_{b \times 1}g = X \bigoplus_{b \times k} j$$
$$= X \bigoplus_{k \in k} b_{k \times 1}g$$

y may look like a complicated t-dimensional phenomenon; but in reality it's just k-dimensional. This is understanding! As we saw in the slide 6:

 $X\beta$ is a linear combination of the columns of X. This is a *subspace* of R^t of at most *k* dimensions. The model states that the *t* observations must fall within this subspace.

In other words, y is operating under k, rather than t, degrees of freedom. We have deepened our understanding, if k < t.

When the right β is found, prediction for new Xs is possible.

But reality is usually messy. Models are approximate:

$$y_{b \times 1}g \approx f \mathbf{\Phi}_{b \times 1}g$$

The equality is regained by adding a random error term; so the model becomes a *statistical* model:

$$\mathbf{y}_{\mathbf{b}\times 1}\mathbf{g} = f \mathbf{e}_{\mathbf{b}\times 1}\mathbf{g}\mathbf{j} + \mathbf{e}_{\mathbf{b}\times 1}\mathbf{g}$$

Specifically, a *linear* model becomes a *linear* statistical model:

$$\mathbf{y}_{\mathbf{b}\times 1}\mathbf{g} = \mathbf{X}_{\mathbf{b}\times k}\mathbf{g}\mathbf{b}_{\mathbf{b}\times 1}\mathbf{g}^{\mathbf{+}\mathbf{e}}\mathbf{b}_{\mathbf{b}\times 1}\mathbf{g}$$
$$(E[\mathbf{e}] = 0. \text{ And let } Var[\mathbf{e}] = \Sigma = \sigma^2 \Phi, \text{ which is symmetric } t \times t.)$$

Two in-depth papers by the author on estimating the β of the linear statistical model and on predicting:

1. "Loss Prediction by Generalized Least Squares," *PCAS* LXXXIII (1996), 436-489.

2. "Conjoint Prediction of Paid and Incurred Losses," CAS Forum, Summer 1997, 241-379.

The author's "Bible" on the subject:

George G. Judge, *et al.*, Introduction to the Theory and Practice of Econometrics, 2nd edition (Wiley, 1988).

But here follows a "quick and dirty" derivation of the estimator ...

$$\mathbf{y}_{\mathbf{b}\times 1}\mathbf{g} = \mathbf{X}_{\mathbf{b}\times k}\mathbf{g}\mathbf{b}_{\mathbf{b}\times 1}\mathbf{g}^{\mathbf{+}}\mathbf{e}_{\mathbf{b}\times 1}\mathbf{g}$$
$$\mathbf{W}_{\mathbf{y}} = \mathbf{W}_{\mathbf{b}\times t}\mathbf{g}\mathbf{X}_{\mathbf{b}\times k}\mathbf{g}\mathbf{b} + \mathbf{W}_{\mathbf{e}}$$

If the square $(k \times k)$ matrix WX has an inverse (non-singular):

$$b_{VX}g_{W}y = b_{VX}g_{WX}b + b_{VX}g_{W}e$$

$$b_{VX}g_{W}y = b + b_{VX}g_{W}e$$

$$E[b_{VX}g_{W}y] = E[b + b_{VX}g_{W}e]$$

$$= b$$

The last equality holds because β , W, and X are constants, and $E[\mathbf{e}] = 0$.

So we have a *L*inear-in-y and *U*nbiased *E*stimator of β :

$$\hat{\beta} = \mathbf{b} \mathbf{V} \mathbf{X} \mathbf{g} \mathbf{W} \mathbf{y}$$

$$= \mathbf{b} \mathbf{V} \mathbf{X} \mathbf{g} \mathbf{W} \mathbf{b} \mathbf{x} \mathbf{b} + \mathbf{e} \mathbf{g}$$

$$= \mathbf{b} \mathbf{V} \mathbf{X} \mathbf{g} \mathbf{W} \mathbf{x} \mathbf{b} + \mathbf{b} \mathbf{V} \mathbf{X} \mathbf{g} \mathbf{W} \mathbf{e}$$

$$= \mathbf{b} + \mathbf{b} \mathbf{V} \mathbf{X} \mathbf{g} \mathbf{W} \mathbf{e}$$

$$E[\hat{\beta}] = \mathbf{b}$$

$$Var[\hat{\beta}] = Var[\mathbf{b} \mathbf{V} \mathbf{X} \mathbf{g} \mathbf{W} \mathbf{e}]$$

In the references it is derived that Var[Ae] = AVar[e]A'. Hence:

$$Var[\hat{\beta}] = Var[\mathbf{b} X \mathbf{g} W \mathbf{e}]$$

= $\mathbf{b} X \mathbf{g} W Var[\mathbf{e}] W' \mathbf{b} Y' W' \mathbf{g}$
= $\mathbf{b} X \mathbf{g} W \Sigma W' \mathbf{b} Y' W' \mathbf{g}$

In the special case that $W\Sigma W' = X'W'$ (or $W = X'\Sigma^{-1}$):

$$Var[\hat{\beta}] = \mathbf{b} \mathbf{x} \mathbf{g}^{\dagger}$$
$$= \mathbf{c} \mathbf{x}' \Sigma^{-1} \mathbf{x} \mathbf{h}^{\dagger}$$

This special W exploits the variance of **e**, so that:

$$\forall W \mathbf{x}' \Sigma^{-1} X \mathbf{h} \leq \mathbf{b} X \mathbf{g} W \Sigma W' \mathbf{b} Y' W' \mathbf{g}$$

The inequality is meaningful in the context of non-negative definite matrices (Appendix A of paper).

The special W makes for the **B** est **L** inear **U** nbiased **E** stimator:

$$\hat{\boldsymbol{\beta}} = \boldsymbol{b} \boldsymbol{W} \boldsymbol{X} \boldsymbol{g}^{\dagger} \boldsymbol{W} \boldsymbol{y}$$
$$= \boldsymbol{C} \boldsymbol{X}' \boldsymbol{\Sigma}^{-1} \boldsymbol{X} \boldsymbol{h}^{\dagger} \boldsymbol{X}' \boldsymbol{\Sigma}^{-1} \boldsymbol{y}$$
$$= Var [\hat{\boldsymbol{\beta}}] \boldsymbol{X}' \boldsymbol{\Sigma}^{-1} \boldsymbol{y}$$

The estimator is invariant to the scale of W. For any scalar $\alpha \neq 0$:

$$b = b X g^{-1} a W y = b X g^{$$

So the **BLUE** of β is invariant to the scale of *Var*[**e**]:

$$\hat{\boldsymbol{\beta}} = \mathbf{C} \mathbf{X}' \boldsymbol{\Sigma}^{-1} \mathbf{X} \mathbf{h}^{1} \mathbf{X}' \boldsymbol{\Sigma}^{-1} \mathbf{y}$$

$$= \mathbf{C} \mathbf{X}' \mathbf{G}^{2} \boldsymbol{\Phi} \mathbf{h}^{1} \mathbf{X} \mathbf{j}^{-1} \mathbf{X}' \mathbf{G}^{2} \boldsymbol{\Phi} \mathbf{h}^{1} \mathbf{y}$$

$$= \mathbf{C} \mathbf{X}' \boldsymbol{\Phi}^{-1} \mathbf{X} \mathbf{h}^{1} \mathbf{X}' \boldsymbol{\Phi}^{-1} \mathbf{y}$$

Later on σ^2 can be estimated (page 69 of paper). Shape, not scale!

The simplest shape is the identity, or $\Phi = I_t$. Then the BLUE is:

$$\hat{\boldsymbol{\beta}} = \mathbf{C} \mathbf{X}' \mathbf{I}_{t}^{-1} \mathbf{X} \mathbf{h} \mathbf{X}' \mathbf{I}_{t}^{-1} \mathbf{y} = \mathbf{b} \mathbf{X}' \mathbf{X} \mathbf{g} \mathbf{X}' \mathbf{y}$$

To use the identity shape when the true shape is otherwise — perhaps we're ignoring information — is to settle for a not-as-good estimator; but it's still an unbiased estimator.

This "quick and dirty" approach (with the $k \times t$ matrix W) is related to the interesting subject of instrumental variables. Consult Judge's Econometrics textbook, pages 577-579.

Now follows a geometrical interpretation of the linear statistical model, otherwise known as *least squares* ...

$$\hat{\beta} = \mathbf{b} \mathbf{x}' \mathbf{X} \mathbf{g} \mathbf{X}' \mathbf{y}$$

can be shown (Judge, 190-192) to minimize the function:

which function represents the square of the Euclidean distance between **y** from X β . But X_(t×k) β _(k×1) is a *k*-dimensional subspace of the *R^t* which **y** inhabits. (Recall from slide 18 that X β is a linear combination of the columns of X.)

To minimize $f(\beta)$ is to find the point closest to y of the subspace spanned by X. At this point, at this particular X β , y drops a perpendicular to the subspace: $(\mathbf{Y} - X\hat{\beta}) \perp X\hat{\beta}$ The closer **y** is to the subspace, the more tempted we are to say that **y** *is* $X \hat{\beta}$ (barring a little randomness, which we can quantify and manage). And a *t*-dimensional phenomenon is more or less reduced to *k* dimensions. To repeat, this is a deeper understanding, if *k* < *t*.

In general,
$$\hat{\beta} = \mathbf{X}' \Phi^{-1} \mathbf{X} \mathbf{h} \mathbf{X}' \Phi^{-1} \mathbf{y}$$
 minimizes:
 $f \mathbf{b} \mathbf{g} = \mathbf{b} - \mathbf{X} \mathbf{b} \mathbf{g} \Phi^{-1} \mathbf{b} - \mathbf{X} \mathbf{b} \mathbf{g}$

This represents a generalized Euclidean distance between \mathbf{y} and $X\beta$, since Φ will penalize differences in some directions more heavily than differences in other directions. With a non-identity (and positive definite) Φ , constant distance from a center takes the form of an ellipse, rather than that of a circle.

Prediction

Estimating the parameter (β) of a statistical model usually isn't enough. Typically, we'll want to predict new **y**s, given new Xs.

Also, we ought to know how much the phenomenon can vary from our prediction — to know the variance of the prediction from its expected value.

Predictions can be correlated with what we've observed. We can't always have the simple world of i.i.d. (independent, identically distributed)

The formulation, with the help of partitioned matrixes:

 $_1$ rows of observations, t_2

1 is t_1 is $t_1 \times \beta$ is $k \times 1$, is $t_1 \times \phi_{12}$ is $\times t_2$, etc. The Φ matrix is symmetric, so Φ_{12} 21'.

 \mathbf{y}_1 s are known (or taken for granted).

 \mathbf{y}_2 contains missing values. We want to estimate (or predict) it.

$$\mathbf{y}_2$$
 is ...

$$\hat{\mathbf{y}}_{2} = X_{2}\hat{\boldsymbol{\beta}} + \Phi_{21}\Phi_{11}^{-1}\boldsymbol{\mathfrak{G}}_{1} - X_{1}\hat{\boldsymbol{\beta}}\boldsymbol{j}$$
where $\hat{\boldsymbol{\beta}} = \boldsymbol{\mathfrak{X}}_{1}^{\prime}\Phi_{11}^{-1}X_{1}\boldsymbol{h}^{\dagger}X_{1}^{\prime}\Phi_{11}^{-1}\mathbf{y}_{1}$

$$Var[\mathbf{y}_{2} - \hat{\mathbf{y}}_{2}] = \boldsymbol{s}^{2}\boldsymbol{\mathfrak{G}}_{22} - \Phi_{21}\Phi_{11}^{-1}\Phi_{12}\boldsymbol{h}$$

$$+\boldsymbol{\mathfrak{X}}_{2} - \Phi_{21}\Phi_{11}^{-1}X_{1}\boldsymbol{h}ar[\hat{\boldsymbol{\beta}}]\boldsymbol{\mathfrak{K}}_{2} - \Phi_{21}\Phi_{11}^{-1}X_{1}\boldsymbol{h}$$
where $Var[\hat{\boldsymbol{\beta}}] = \boldsymbol{s}^{2}\boldsymbol{\mathfrak{K}}_{1}^{\prime}\Phi_{11}^{-1}X_{1}\boldsymbol{h}^{\dagger}$

 $\Phi_{21} \neq 0$ allows errors in the observations to affect the predictions. Looks nasty, but really quite gentle. Proof in Appendix C of "Conjoint Prediction." Also see pages 68f. of this paper.

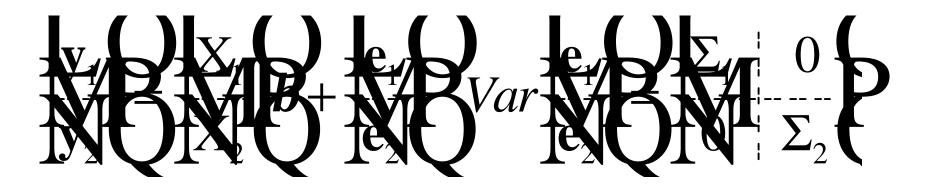
3. Credibility and Prior Information

"Readers who have come this far may conclude from what they've read that casualty actuarial science is the study and application of the theory of credibility, and that's all. Is it all?"

Matthew Rodermund, Foundations of Casualty Actuarial Science, 19.

It's hard to answer "No" to Rodermund's question. Actuaries love the $'Z \times A + (1-Z) \times E'$ credibility formula. It blends observation ('A' for 'actual') and prior opinion ('E' for 'expected').

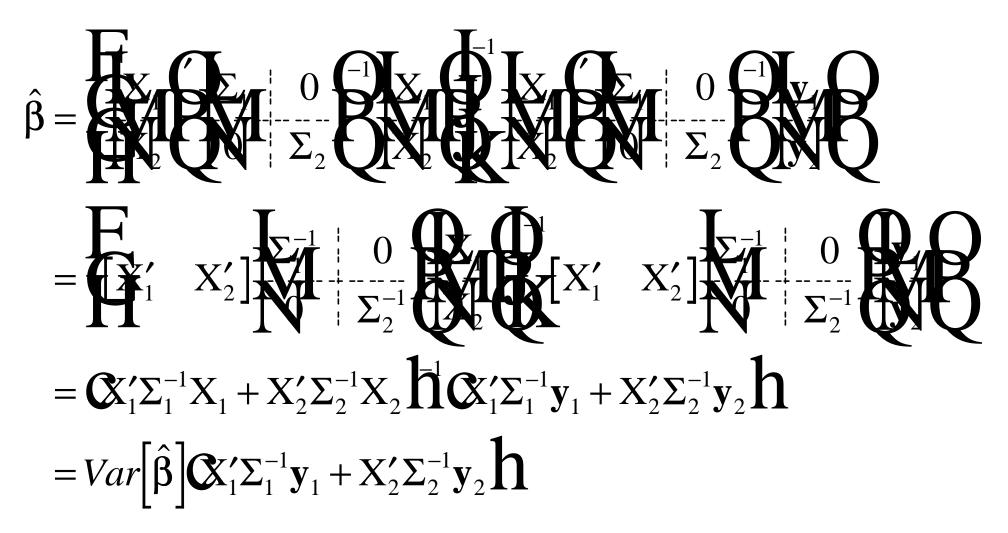
But a (linear) statistical model does the blending even better ...



Observations only: t_1 rows from source 1, and t_2 rows from source 2.

Simple variance structure in that the sources do not covary (off-diagonal Σ s are zero).

The BLUE of β is ...



Not too bad when one has a feel for partitioned matrices (slide 14).

Next to last equation looks like a (matrix) weighted average. This can be made explicit ...

$$\hat{\boldsymbol{\beta}} = \mathbf{X}_{1}'\Sigma_{1}^{-1}X_{1} + X_{2}'\Sigma_{2}^{-1}X_{2}\mathbf{h} \mathbf{X}_{1}'\Sigma_{1}^{-1}\mathbf{y}_{1} + X_{2}'\Sigma_{2}^{-1}\mathbf{y}_{2}\mathbf{h}$$

$$= \mathbf{X}_{1}'\Sigma_{1}^{-1}X_{1} + X_{2}'\Sigma_{2}^{-1}X_{2}\mathbf{h} \mathbf{X}_{1}'\Sigma_{1}^{-1}X_{1}\mathbf{X}_{1}\mathbf{X}_{1}'\Sigma_{1}^{-1}X_{1}\mathbf{h}X_{1}'\Sigma_{1}^{-1}\mathbf{y}_{1} + \cdots \mathbf{j}$$

$$= \mathbf{X}_{1}'\Sigma_{1}^{-1}X_{1} + X_{2}'\Sigma_{2}^{-1}X_{2}\mathbf{h} \mathbf{X}_{1}'\Sigma_{1}^{-1}X_{1}\mathbf{\beta}_{1} + X_{2}'\Sigma_{2}^{-1}X_{2}\mathbf{\beta}_{2}\mathbf{j}$$

$$= \mathbf{X}_{1}'\Sigma_{1}^{-1}X_{1} + X_{2}'\Sigma_{2}^{-1}X_{2}\mathbf{h} \mathbf{X}_{1}'\Sigma_{1}^{-1}X_{1}\mathbf{\beta}_{1} + X_{2}'\Sigma_{2}^{-1}X_{2}\mathbf{\beta}_{2}\mathbf{j}$$

The estimator of the two-source model weights the estimators of the one-source models according to the inverses of the variances of those estimators (harmonic average — Appendix A).

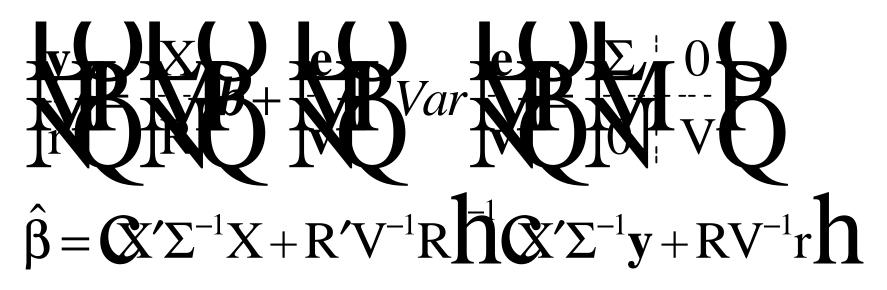
Since β is $k \times 1$, this is credibility in *k* dimensions.

$$Var[\hat{\beta}] = \bigoplus ar^{-1}[\hat{\beta}_{1}] + Var^{-1}[\hat{\beta}_{2}]\mathbf{j}^{-1}$$
$$Var^{-1}[\hat{\beta}] = Var^{-1}[\hat{\beta}_{1}] + Var^{-1}[\hat{\beta}_{2}]$$
$$> Var^{-1}[\hat{\beta}_{1}], Var^{-1}[\hat{\beta}_{2}]$$
$$Var[\hat{\beta}] < Var[\hat{\beta}_{1}], Var[\hat{\beta}_{2}]$$

As Appendix A explains (using positive definite matrices), the twosource estimator is of less variance than either of the one-source estimators. The more knowledge, the better.

The second source doesn't actually have to be observed. It can be theory, opinion, or guess — anything on which you're willing to rely.

The second source, or the prior information, doesn't even have to be a complete model, so that a β_2 could be estimated:



If $k \times k \operatorname{R'V}^{-1} \operatorname{R}$ is of rank $j \le k$, the variance of the two-source estimator will be improved along j orthogonal axes. If j < k, then $\operatorname{R'V}^{-1} \operatorname{R}$ is singular and the second source cannot produce its own estimate for β . But it still improves the mixed estimator in j out of k dimensions. See Appendix A. Sorry, no example is given in this presentation. But the paper works out several examples.

Benefits of statistical modeling to credibility:

- Provides an systematic and orderly framework.
- Furnishes variances of estimates and predictions, as well as means.
- Extends credibility from one to *k* dimensions.

4. Credibility and Random Effects

Hardest part of the paper — Sections 9 and 10, and Appendix E. Given n related groups, with non-covarying **e**s and **v**s:

$$\mathbf{y}_{i} \mathbf{b}_{\times 1} \mathbf{g} = \mathbf{X}_{i} \mathbf{b}_{\times k} \mathbf{g}^{k} \mathbf{b}_{i} \mathbf{b}_{\times 1} \mathbf{g}^{k} \mathbf{e}_{i} \mathbf{b}_{\times 1} \mathbf{g}^{k} \operatorname{Var}[\mathbf{e}_{i}] = \Sigma_{i} \mathbf{b}_{\times k} \mathbf{g}$$

But $\mathbf{b}_{i} = \mathbf{b}_{0} + \mathbf{v}_{i}$, $\operatorname{Var}[\mathbf{v}_{i}] = V_{\mathbf{b} \times k} \mathbf{g}$
So $\mathbf{y}_{i} = \mathbf{X}_{i} \mathbf{b}_{0} + \mathbf{b}_{i} \mathbf{v}_{i} + \mathbf{e}_{i} \mathbf{g}$
 $= \mathbf{X}_{i} \mathbf{b}_{0} + \mathbf{\tau}_{i}$, $\operatorname{Var}[\mathbf{\tau}_{i}] = \mathbf{X}_{i} \operatorname{VX}_{i}' + \Sigma_{i}$

This is just a linear statistical model, so we can estimate β_0 , the β_i s, and any predictions built on the β_i s.

If V is unknown, it may be estimated by the method of variance components (ML method also possible — Appendix F).

If V is large, the β_i s are free and the groups have much credibility. If V is small, the β_i s are close to β_0 and the groups have little credibility.

Appendix E discusses the random-effects model in detail. A beautiful result is that the simple average of the estimates of the β_i s must equal the estimate of β_0 . In effect, credibility democratizes the groups.

Section 10 presents a random-effects trend model, a twodimensional credibility problem.

5. Conclusion

Arthur Bailey challenged (1945-1950) classical statistics with three problems which justified his "greatest accuracy credibility:"

- Use of prior information in estimation $(Z \times A + (1-Z) \times E)$
- Estimating for an individual that belongs to a heterogeneous population (merit and experience rating a fruitful subject for Bayesian credibility. See Appendix B)
- Estimating for groups together, which is more accurate than estimating each separately (the random-effects model. See Sections 9 and 10, and Appendix E)

The paper shows how modern statistics solves these problems, to the legitimization and enrichment of credibility.

Corrections

Page:Line	Text	Change
74:14	0.904	0.113
74:15	4.240	0.606
87:1	ends of the introduction and	end
118:3	element	variance
121:1	subscript <i>i</i> is ranges	subscript <i>i</i> ranges
137:2	$\sum_{i=1}^{n} ((\mathbf{y}_{i} - \hat{\boldsymbol{\mu}}) - (\boldsymbol{\nu} - \hat{\boldsymbol{\mu}}))((\mathbf{y}_{i} - \hat{\boldsymbol{\mu}}) - (\boldsymbol{\nu} - \hat{\boldsymbol{\mu}}))$	$\sum_{i=1}^{n} \left(\left(\mathbf{y}_{i} - \hat{\boldsymbol{\mu}} \right) - \left(\boldsymbol{\nu} - \hat{\boldsymbol{\mu}} \right) \right) \left(\left(\mathbf{y}_{i} - \hat{\boldsymbol{\mu}} \right) - \left(\boldsymbol{\nu} - \hat{\boldsymbol{\mu}} \right) \right)'$
147:11 and 14 (three times)	$X_i T_i^{-1} X_i'$	$\mathbf{X}_{i}^{\prime}\mathbf{T}_{i}^{-1}\mathbf{X}_{i}$

Vice President of Actuarial Research and Development American Re-insurance Company (as of March 16, lhalliwell@amre.com)