# The Philosophical Significance of Stein's Paradox

Olav Vassend\*, Elliott Sober\*, and Branden Fitelson<sup>+</sup>

\*Philosophy Department, University of Wisconsin, Madison, Wisconsin †Philosophy Department, Northeastern University, Boston, Massachusetts

*Abstract*: Charles Stein discovered a paradox in 1955 that many statisticians think is of fundamental importance. Here we explore its philosophical implications. We outline the nature of Stein's result and of subsequent work on shrinkage estimators; then we describe how these results are related to Bayesianism and to model selection criteria like the Akaike Information Criterion. We also discuss their bearing on scientific realism and instrumentalism. We argue that results concerning shrinkage estimators underwrite a surprising form of holistic pragmatism.

# **1.** Shrinkage is better than straight MLE when $k \ge 3^1$

If you sample at random (with replacement) from a human population and find that the average height in your sample is 5 feet, what could be more natural than the conclusion that the average height in the whole population is about 5 feet? The principle underwriting this inference has gone by different names. Philosophers have called it "the principle of induction." Frequentist statisticians say that the inference is justified by a method called "maximum likelihood estimation" (MLE). Here the word "likelihood" is used in its technical sense. The estimate that the population mean is 5 feet maximizes likelihood, not in the sense that this is the most probable estimate given the observations, but in the sense that it makes the observations more probable than other estimates are able to do. The likelihood of hypothesis H relative to observation O is the quantity Pr(O|H), not the quantity Pr(H|O). If the population mean were 7 feet, a sample mean of (about) 5 feet would be very improbable, "almost a miracle." If the population mean were 5 feet, a sample mean of approximately 5 feet would be much less surprising.<sup>2</sup>

<sup>&</sup>lt;sup>1</sup> Here we are indebted to the excellent exposition in Efron and Morris (1977).

<sup>&</sup>lt;sup>2</sup> Here we assuming the variances are the same in the two cases.

Statisticians have done for MLE something that their philosophical predecessors did not do for the principle of induction. They proved that MLE uniquely possesses various desirable properties. Gauss showed that if the distribution of heights in the population is normal, then the ML estimate of the mean height is the sample mean (Edwards 1974, p. 11). Gauss also realized that MLE is unbiased, meaning that repeated ML estimates based on different samples drawn from the same population will tend to be centered on the population mean. There are infinitely many unbiased estimators; Gauss (1823) proved, finally, that ML estimates of normal means have lower expected mean-squared error than any other unbiased estimator that is a linear function of the observations.

The case for MLE was strengthened in the 1930's with the development of statistical decision theory. Suppose estimates are good to the degree that they come close to the true (but unknown) value of the quantity being estimated. In particular, consider a "loss function" that measures an estimator's inaccuracy by its squared deviation from the true value of the target quantity. Figure 1 depicts three estimators that might be used for the mean height in a population - the sample average, one-half the sample average, and the sample median (the middle value). Figure 1 tracks how accurate an estimator can be expected to be as a function of the unknown value of the population's mean height ( $\theta$ ). Of course, how well an estimator does may vary from sample to sample; what can be plotted precisely is the average performance (the mathematical expectation). The greater the expected squared error, the higher the estimator's "risk." As Figure 1 makes plain, the mean has a lower expected inaccuracy than the median for each possible value of  $\theta$ . The mean *strongly dominates* the median, in the technical sense of that term used in decision theory. For that reason, the median is deemed an *inadmissible* estimator. An admissible estimator isn't dominated (either strongly or weakly) by any other estimator.<sup>3</sup> Figure 1 does not say whether the sample mean and half-the-sample-mean are admissible, since the figure leaves open that there might be other estimators that dominate both. Blyth (1951) and Lehmann and Hodges (1951) settled this question; they proved, for data drawn from a normal

<sup>&</sup>lt;sup>3</sup> X weakly dominates Y precisely when X's risk is never higher than Y's, and for some values of  $\theta$ , X's risk is lower. Inadmissible estimators are weakly dominated; they may or may not be strongly dominated.

population, that there is no estimator that dominates MLE. MLE is admissible. This doesn't mean that MLE dominates all other estimators, but only that no estimator dominates MLE.





Stein (1956) surprised the statistics community by showing that MLE has a dark side. He proved that MLE is inadmissible when the estimation problem concerns three or more statistically independent measurements of the means of normal distributions with identical known variances. MLE is admissible if you are estimating average height in the United States. It also is admissible if you are estimating average height in Norway. And it is admissible if you are estimating average height in lapan. But if you want to estimate all three averages at the same time, and your goal is to minimize expected error across the three estimates, MLE is inadmissible.<sup>5</sup> James and Stein (1961) showed that you do better in expectation if you construct

<sup>&</sup>lt;sup>4</sup> Adapted from Efron and Morris (1977), p. 124.

<sup>&</sup>lt;sup>5</sup> In the theory of random walks, a similarly interesting transition happens when you move from two to three dimensions. Random walks in one or two dimensions are recurrent, meaning they have a probability of 1 of returning to their starting point. Random walks in dimensions greater

your estimates by shrinking the three observed frequencies towards zero by multiplying each by  $c = 1 - \sigma^2 / (\sum_{i=1}^3 X_i^2)$ . Here,  $X_i$  is the ML estimate of the average height of population *i*, and  $\sigma^2$  is the variance within each population. Multiplying by *c* shrinks estimates towards 0 because *c* will, in general, be less than 1 and greater than 0.<sup>6</sup> Efron and Morris (1973) showed how to generalize Stein's result when the variances of the populations are unknown and/or different from each other, and they also showed that shrinking towards the grand mean of the samples is better than straight MLE when four or more quantities are estimated.<sup>7</sup>

When you estimate three or more parameters, you should not expect a shrinkage estimator to better MLE for *each* parameter; in fact, for each parameter, there are parameter values for which shrinkage will do worse than MLE and other values for which it will do better. However, your *total* error over the parameters *collectively* will be lower in expectation if you shrink.

# 2. Why this is paradoxical

Stein's result is bizarre. The three quantities can be entirely unrelated to and independent of each other and these shrinkage estimators still do better than straight MLE when judged by the criterion of expected squared error. Efron and Morris (1977) give the example of estimating the batting abilities of 18 Major League baseball players from data on their first 45 times at bat in a year. If you estimate each player's probability of getting a hit, and use those estimates to predict how well each player will do at the end of the year, the data from the end of the season reveals that you do worse if you use MLE than if you shrink those ML estimates towards the grand sample mean. In this case, real data behave in conformity with Stein's result. This may lead you to think that the baseball players are influenced by a common cause that exerts a "gravitational

than two are not recurrent. Brown (1971) discovered a connection between Stein's result and this fact about random walks.

<sup>&</sup>lt;sup>6</sup> Although it's unlikely, it's possible for c to be a negative number.

<sup>&</sup>lt;sup>7</sup> The estimator that shrinks the ML estimates towards the sample mean is now standardly referred to as the "Efron-Morris estimator," and we will follow this practice in our paper; however, it is worth pointing out that the Efron-Morris estimator was originally suggested by Dennis Lindley in the discussion section of Stein (1962).

attraction" on batting ability as the season unfolds. The point of importance is that shrinkage results depend on no such assumption (Stigler 1990, pp. 147-148).

Efron and Morris (1977) make this point graphic. After describing their baseball example, they add to it the problem of inferring the percentage of foreign cars in Chicago. This enlarged problem involves quantities that come in different units  $-\frac{\text{number of hits}}{\text{number of times at bat}}$  and  $\frac{\text{number of foreign cars in Chicago}}{\text{number of cars in Chicago}}$ . Shrinkage estimators are better than straight MLE estimates here, since there are more than three parameters being estimated. The mathematical results that ground this fact do not turn on whether hits in baseball and foreign cars in Chicago are causally related to each other.

|                         | Which quantity do you want to estimate? |                   |
|-------------------------|---|-------------------|
|                         | % domestic cars                         | % foreign cars    |
| Your data               | 60% domestic cars                       | 40% foreign cars  |
| Your shrinkage estimate | <60% domestic cars                      | <40% foreign cars |

This is weird, but there is more. The additional strangeness is that there are different, equally correct, ways of coding your data, and your choice of code will affect how you shrink your estimates away from the MLE estimates. This point is illustrated in the accompanying table. The average success rate at the start of the season for the 18 baseball players that Efron and Morris examined was 26.5%. Suppose you sample the cars in Chicago and find that 40% of the cars in your sample are foreign. Since there are 18 baseball players and only one Chicago, the grand mean of these 19 frequencies is close to 27%. So, if you shrink your estimates of those 19 parameters towards that grand mean, you'll do better in expectation than if you use straight MLE. Notice that if you shrink, you will shrink your car estimate towards 27% regardless of whether you estimate the percentage of foreign cars or the percentage of domestic. If you do both, you'll have contradictory estimates. So what should you do? Of course, only one of these shrinkages will move you closer to the truth, but both in expectation will do so when they are part of the 19-parameter problem. Moreover, what is true in the 19-parameter problem also holds for the initial problem of the 18 baseball players. The James-Stein theorem does not require that you code the 18 sample means in terms of each player's percentage of hits in their first 45 times at bat. You could code some players in terms of their success rates and others in

terms of their failure rates, and the theorem would still apply. As Efron and Morris tell the story, Roberto Clemente's estimated season-long batting average gets shrunk from his initial success rate of 0.400 towards 0.265, but MLE also will be bettered if you shrink his initial failure rate of 0.600 towards 0.265. What goes for Chicago also goes for Roberto.

#### 3. Stein's Paradox and Linguistic Invariance

The reader has a right to be shocked by the fact that what you say about Roberto Clemente's batting ability depends on whether you seek to estimate his ability to succeed or his tendency to fail. The oddity arises if your problem is to estimate Clemente's ability and that of 17 other baseball players simultaneously. This finding shows that shrinkage estimators are not linguistically invariant in the following weak sense:

Linear Invariance: An estimator E of parameter  $\theta$  is linearly invariant iff, for any data set D, and any linear function f,  $f{E[\theta | D]} = E[f(\theta) | D$ , where a linear function is any function of the form  $f(\theta) = a\theta+b$ , for real-valued a and b.

Linear invariance seems like a very reasonable requirement to impose on estimators. Demanding linear invariance means that if your estimator says that 5 feet is the best estimate of the mean height in a population, it had better say that 60 inches is also the best estimate. MLE is invariant in this sense, as are the sample mean, the sample median, and the sample mode.

There is another sort of invariance that also seems reasonable to impose on estimators. Suppose you are estimating the mass m of some object and that you end up with a data set D of measurements. Now suppose someone adds the number 1 to each of your measurements, thereby yielding the new data set D'. Intuitively, it seems reasonable in this case to demand that your estimator should be invariant in the following sense: E(m | D') = E(m | D) + 1. This type of invariance is known as *translation invariance*. Note that whereas Linear Invariance is invariance under certain transformations of the parameters, translation invariance is invariance under certain transformations of the sample mean, median, and mode are all clearly translation invariant, and MLE satisfies translation invariance too given very plausible restrictions on the distribution. Shrinkage estimators, however, are not translation invariant.

Given weak restrictions on the loss function and the probability distribution, any translation invariant estimator has a constant risk function (cf. the median and the average in Figure 1). Translation invariant estimators are therefore easy to compare with each other: the *best* translation invariant estimator is simply the one that has the lowest (constant) risk. Stein's result shows that, under quadratic loss, the best translation invariant estimator of normal means, namely MLE, is dominated by shrinkage estimators that are not translation invariant and not invariant under linear transformations, when more than two parameters are being estimated.

MLE is the best invariant estimator if you are estimating the mean of a normal distribution and you are using the squared loss function to measure inaccuracy. However, if you change your loss function, then MLE might no longer be the best invariant estimator. For example, if you were to measure inaccuracy by taking the absolute value of the deviation from the target quantity, rather than the square of this deviation, then the best invariant estimator would instead be the sample median, while the ML estimator would still be the sample mean. On the other hand, if you were measuring the mean of a Laplace distribution rather than a normal distribution, then the ML and best estimator would be the median.<sup>8</sup>

The reader might naturally wonder whether Stein's result is a mathematical curiosity that arises only in estimation problems that involve normal distributions where the loss function is squared deviation. It is not. Brown (1966) showed that for a very large range of loss functions and distributions, the best translation invariant estimator is *never* admissible. In each of the cases considered by Brown, there is a way of modifying the James-Stein estimator to yield a non-invariant estimator that dominates the best invariant estimator.<sup>9</sup>

<sup>&</sup>lt;sup>8</sup> The Laplace distribution is exponential like the normal distribution. The main difference between the two distributions is that in the Laplace distribution, the exponent is not squared; instead, its absolute value is taken.

<sup>&</sup>lt;sup>9</sup> Nor is Stein's result in the case of the squared loss function attributable to the fact that the squared loss function is unbounded (James and Stein 1961, p. 367). The result also does not depend on the assumption that the measurements of the different means are independent of one another (Bock 1975).

Returning to the baseball players, it seems natural to restrict your attention to the class C of estimation methods that yield the same answer whether you code your data in terms of failure rates or success rates, and that aren't affected by uniformly shifting the data by some constant amount. However, Brown's result shows that, under widely applicable conditions, the best estimators in C will be dominated by strange shrinkage estimators that are sensitive to how you encode your data.

The real Stein paradox is therefore not just that MLE is inadmissible in dimensions greater than 2. The real paradox is that if total accuracy is your goal, then under widely applicable conditions, you have to go outside the class of invariant estimators if you want your estimator to be admissible.<sup>10</sup>

If you value linguistic invariance above total accuracy, you might choose to use an invariant estimator even though shrinkage is more accurate.<sup>11</sup> However, if you opt for a shrinkage estimator SE(-), there is a kind of invariance that your estimator will possess, at least if you use the squared loss function. If SE( $\theta$ , data) is more accurate in expectation than MLE( $\theta$ , data), then SE( $\theta$ ', data) will also be more accurate in expectation than MLE( $\theta$ ', data), provided that  $\theta$ ' is a linear transformation of  $\theta$ . In this sense, it doesn't matter whether you score all 18 baseball players by their success rates, or all by their failure rates, or use success for some and failure for others. Regardless of coding, shrinkage can be expected to better MLE.

# 4. Making the shrinkage result intuitive<sup>12</sup>

Suppose you make independent measurements of three unknown parameters where the measurement of each parameter is modeled as a normal distribution with variance 1 and mean  $\theta_i$ . That is, you model the measurements  $X_i$  by the formula  $X_i = \theta_i + \text{error}$ , where the error is

<sup>&</sup>lt;sup>10</sup> We thank an anonymous reviewer for pointing this out to us.

<sup>&</sup>lt;sup>11</sup> George Barnard seems to endorse this view (see Stein 1962, p. 288), as do Perlman and Chaudhuri (2012, p.139n18), who maintain that shrinkage estimation should not be used in the absence of prior information that non-arbitrarily singles out some point towards which you shrink the ML estimates.

<sup>&</sup>lt;sup>12</sup> Here we are indebted to Stigler (1990).

Gaussian. For simplicity, suppose you have just three measurements,  $x_1$ ,  $x_2$ , and  $x_3$ , of  $\theta_1$ ,  $\theta_2$ , and  $\theta_3$ , respectively. MLE says that you should estimate that  $\theta_i=x_i$  for each i. Figure 2 represents this estimate in the  $\langle X, \theta \rangle$  plane as a 45 degree straight line that goes through the origin. Each  $\hat{\theta}$  is the MLE estimate for a given observed value of X. MLE is therefore a *linear* estimator of  $\theta$  given X, but is MLE the *best* linear estimator?



Figure 2

To investigate this question, it is useful to indulge in a fiction: suppose you know the true value for  $\theta$  that is associated with each of the three observed x values. The three  $\langle x, \theta \rangle$  pairs are represented in Figure 3. What is the best linear estimator given these three data points? To answer this question, you need to decide which of two estimation problems you want to address. These are shown in Figures 3 and 4.

Suppose your goal is to find the line that minimizes the *vertical* distances between points and line. This line is shown in Figure 3; it obeys the equation  $\theta = aX+b$ . Unfortunately, you

don't know the values of the  $\theta$ 's, so you can't estimate the coefficients a and b in this equation in the usual way. However, you can try to approximate them; indeed, all of the different shrinkage estimators, beginning with the one described by James and Stein (1961), can be regarded as clever methods for estimating a and b from the data.



Figure 3

Alternatively, your goal might be to minimize the *horizontal* distances between points and line. The best line is then the one shown in Figure 4. It obeys the equation  $X = \alpha \theta + \beta$ . As Galton (1988) recognized, the lines in Figures 3 and 4 are different. You need to decide whether you want a line that, for an observed x value, is close to the true  $\theta$  value (Figure 3), or a line that, for a given  $\theta$  value, is close to the observed x value (Figure 4).





The ML estimate shown in Figure 2 may be regarded as an approximation of the leastsquares line in Figure 4. Indeed, the 45 degree ML estimate is equivalent to the theoretical regression line of X on  $\theta$ , and therefore may be said to be the best possible approximation of the least squares line in Figure 4, given that the  $\theta$  values are unknown. However, if you want to minimize error in estimating  $\theta$ , then the line in Figure 4 is not the line you should try to approximate; instead, you should try to approximate the least squares line in Figure 3, which is what shrinkage estimators all attempt to do.

The fact that there are two least-squares lines explains why MLE leaves something to be desired if your goal is to find a straight line that is close to the true  $\theta$  value for a given observed value of X. However, it explains more than that.

First, note that if you were to move the  $\theta$ s in Figure 3 and Figure 4 further apart from each other, the two regression lines would then converge toward each other. Shrinkage can therefore be expected to offer only a minor improvement over MLE if the  $\theta$ s are far apart from each other (although the improvement will always be greater than 0). On the other hand, if the  $\theta$ s are very close together, the slopes of the two regression lines will be very different, and shrinkage estimation can therefore be expected to offer a more dramatic improvement over MLE.

Second, suppose there are just two parameters ( $\theta_1$  and  $\theta_2$ ) that you want to estimate. In that circumstance, the two least squares lines collapse into one, as shown in Figure 5. Since the ML line is the best approximation of the X= $\alpha\theta$ + $\beta$  least squares line, and since this line is necessarily identical to the  $\theta$ =aX+b least squares line, it follows that the ML line must also be the best possible approximation of the  $\theta$ =aX+b least squares line. In other words, the ML estimator is admissible by coincidence (literally) when you try to estimate fewer than three parameters.



Figure 5

We hope this connection of shrinkage estimators to the two regression problems – fitting  $\theta$  to X and fitting X to  $\theta$  – is instructive. However, it is no substitute for the detailed mathematics of James and Stein (1961) and of Efron and Morris (1973), where a particular shrinkage estimator is shown to have lower expected mean squared error than straight MLE. Our point here is to show why MLE is suboptimal and how constraining three or more estimates to be connected to each other can provide an improvement. The assumption that X and  $\theta$  are related to each other by a straight line is such a constraint. This is not to say that shrinkage estimators assume linearity. Rather, the point is that shrinkage towards some single value is a constraint that links distinct estimation problems to each other; by doing so, shrinkage provides an improvement over straight MLE. In Section 9, we provide a different perspective on how introducing a common constraint can improve estimation.

### 5. Why Care About Admissibility?

So far we've sometimes been a bit careless. We've said that Stein's result asserts that shrinkage estimators are "better" than MLE when more than three parameters are estimated. But, of course, Stein's result is a mathematical theorem and doesn't have any implications regarding what's better or worse unless additional assumptions are added. In particular, Stein's result says that there are shrinkage estimators that *dominate* MLE; this result obviously doesn't allow you to conclude that shrinkage estimators are "better" than MLE unless you add the following normative premise:

(D) If estimator E dominates estimator E', then E is better than E'.

Not everyone will accept (D). For example, Bayesians might reject (D) because the more natural Bayesian approach to estimation is to choose the estimator that maximizes expected utility, and it's not clear whether (D) is consistent with such an approach. We postpone a more thorough discussion of the relation between Stein's paradox and Bayesianism until Section 8. Some frequentists will also want to reject (D). For example, Spanos (2016) argues that the criterion of admissibility and the use of loss functions go against frequentism, and he

consequently (implicitly) rejects (D). Spanos's main objection to (D) is that it will sometimes judge a biased or inconsistent estimator to be better than an unbiased or consistent estimator. In particular, as we shall see in Section 9, all shrinkage estimators are biased. Hence, according to (D), shrinkage estimators are better than ML estimators even though shrinkage estimators are biased.

According to Spanos (2016, p. 15), any biased or inconsistent estimator should be immediately disqualified, because the "primary objective" of estimation is to pinpoint the true value of  $\theta$ , and using a biased or inconsistent estimator means abandoning this goal from the outset.

Of course, unbiased and consistent estimators can also fail to pinpoint the true value of  $\theta$ , given the data at hand, but they will tend to converge on the true value as the data set is enlarged. This asymptotic property is nice for an estimator to have, but is it necessary? The answer is far from clear since the estimation problems we face always involve finite data. And for finite data, no estimator is guaranteed or even likely to get you the true value. We think it is reasonable to prefer one estimator over another when the first has a lower expected error than the second, even if the first estimator is biased while the second one isn't. Of course, this requires that "error" be quantified in terms of some loss function. Spanos asks where the extra "information" is supposed to come from (p. 17) that allows you to pick a particular loss function. It may be true that the choice of a particular loss function is to some extent arbitrary, and that several loss functions may reasonably be said to quantify inaccuracy. However, the fact that the choice of a particular loss functions altogether.

Estimators can be compared along many dimensions: they can be biased or unbiased, they can be consistent or inconsistent, they can be variant or invariant, they can be admissible or inadmissible, etc. One of the lessons of Stein's result is that no estimator will be superior to all other estimators along all the dimensions that you might care about. Thus, ultimately, whether one estimator is "better" than another comes down to which properties you care about and how you weigh them against each other. (D) is therefore not an assumption that is unconditionally true; rather, it is a premise that it is reasonable to accept given that you have certain interests.

In particular, if you care about minimizing your error given finite data – as opposed to making sure you pinpoint the true value of  $\theta$  in the infinite long run limit – then (D) is a reasonable premise for you to accept. But, of course, accepting (D) leads you to prefer shrinkage estimators over MLE and biased estimators over unbiased ones.

#### 6. Choosing an Estimator

If your goal is to minimize expected squared error, which estimator should you use? One possible answer is that if you know that estimator  $E_1$  weakly dominates estimator  $E_2$ , you should not use  $E_2$ . We endorse this answer, with one small qualification that we'll note at the end of this section.<sup>13</sup>

How can this negative advice be supplemented with something positive? We begin with two cautions. The first is that admissibility isn't sufficient for using an estimator. Indeed, in any realistic estimation problem there are infinitely many admissible estimators. For example, suppose you are estimating a single parameter and your estimation method is to just guess that  $\theta=0$  regardless of what your data are. If  $\theta=0$ , your estimator has zero risk. Hence, always guessing that  $\theta=0$  is an estimator that isn't weakly dominated by any other estimator that has positive risk when  $\theta=0$ , and any estimator that is not constant is going to have positive risk. More generally, guessing that  $\theta=c$ , regardless of the data, for any constant number c, will be an admissible estimator.<sup>14</sup> Inadmissibility is bad but admissibility isn't anything special.

<sup>&</sup>lt;sup>13</sup> This is not quite to say that inadmissibility suffices for refusing to use an estimator. You may know that E is inadmissible, but not know the identity of an estimator that weakly dominates E. This may lead you to think that E is better than nothing. We take no stand on whether you should use E or simply refuse to make an estimate.

<sup>&</sup>lt;sup>14</sup> The constant estimator E(X)=c has risk 0 for  $\theta=c$ . Any estimator that weakly dominates E(X)=c must also have risk 0 for  $\theta=c$ . But such an estimator will therefore have a variance of 0, which means (given any reasonable error distribution) that the estimator doesn't vary given different data, and hence it must also be a constant estimator. Thus, any estimator that dominates E(X)=c must itself be a constant estimator, but a constant estimator with  $E(X)=d\neq c$  can't dominate E(X)=c. So all constant estimators of the form E(X)=c are admissible.

The second caution is that admissibility isn't necessary for using an estimator. Perhaps you are in a situation of ignorance. You know that estimator E weakly dominates all the other estimators you have considered, but you don't know whether there exists an as yet unknown estimator that dominates E. In this case you are entitled to use E. That entitlement may lapse if you learn more.<sup>15</sup>

It turns out that while MLE is dominated by the James-Stein estimator when three or more parameters are being estimated, James-Stein is itself dominated by other estimators (Baranchik 1964). The same point holds of the Efron-Morris estimator; it isn't admissible, either (Brown 1971). Are there any shrinkage estimators that are admissible and that dominate MLE? The answer (for  $k\geq 5$ ) is yes; there are several (Strawderman 1971). The implication of what we've just said would therefore seem to be that you should never use either the James-Stein or the Efron-Morris estimator since they are both dominated by other known estimators. However, in practice, statisticians care about computational tractability as well as about minimizing global inaccuracy. Since James-Stein and Efron-Morris are both very simple estimators, and since they have risk functions that are numerically close to the known estimators that they are dominated by (Larry Brown, personal communication), you arguably are justified in using James-Stein or Efron-Morris although these estimators are known to be inadmissible.

#### 7. Holistic Pragmatism

Many epistemologists are *evidentialists* – they think that what you believe (or your degree of belief) should be guided by your evidence and by your evidence alone. Evidentialism has its dissenters. Carnap (1950), for example, argues that some propositions can be accepted because they represent convenient conventions even though there is no evidence that they are true. Pascal (1662) anticipated this pragmatic turn by arguing that belief in God should be

<sup>&</sup>lt;sup>15</sup> Suppose estimators  $E_1$  and  $E_2$  each weakly dominate the others you have considered, but neither dominates the other. You may have reason to prefer one over the other if you have reason to believe that some values of  $\theta$  are more probable than others. Consider the relation of the mean and half-the-mean in Figure 1.

influenced by the positive utility of going to heaven and the negative utility of going to hell. So did James (1896), except that for him the question need not involve the existence of God or the afterlife. Another departure from strict evidentialism may be found in Rudner's (1953) argument that "the scientist *qua* scientist makes value judgments." Rudner maintains that science is in the business of accepting and rejecting hypotheses and that your standards concerning how much evidence is required for you to accept or reject should depend on the ethical consequences of error. None of these pragmatisms covers what Stein discovered. We call this Steinian pragmatism *Holistic Pragmatism* (not to be confused with Morton White's (2005) pragmatism of the same name). Holistic pragmatism is the thesis that when an estimation problem has several parts, it's a pragmatic decision whether your goal is to minimize error across the whole problem, or to minimize error within each part.

Of course, considered separately, MLE and shrinkage estimation are each "evidentialist" in the sense that ML estimators and shrinkage estimators are functions of the evidence and only of the evidence. The reason we think Stein's result licenses a kind of pragmatism is that the result tells you that your goals are relevant to deciding whether you should use MLE or a shrinkage estimator.

Compare Ms. Multi-Tasker and Mr. One-at-a-Timer. Ms. Multi-Tasker takes up three estimation problems and wants to minimize her expected sum of squared errors across the three. Mr. One-at-a-Timer takes up the same three problems, and uses the same evidence and background information that Ms. Multi-Tasker has at hand, but he cares about each problem for its own sake, wanting to minimize his expected squared error on each. According to Stein, they should reach different estimates, with Ms. Multi-Tasker shrinking and Mr. One-at-a-Timer doing straight MLE.

To an evidentialist, whether estimator E is a good estimator of  $\theta$  given data D is purely a question of the relation between E,  $\theta$  and D. Estimation is in other words a three-place relation. However, according to holistic pragmatism, whether your goal is to maximize local or global accuracy also matters. Thus, on holistic pragmatism, estimation involves a four-place relation between E,  $\theta$ , D, and the goals of the agent, G. It is utterly familiar that rational action requires assumptions about utilities. It is controversial that rational belief must involve such assumptions.<sup>16</sup> Some of the utilities involved in deciding whether to use a shrinkage estimator are epistemic – the goal is to have one's estimates be close to the truth (where this is quantified by using the expected sum of the squared errors). Here we see a departure from Pascal (and from James). But an additional type of utility is relevant to estimation: should you care about estimation problems separately or should you seek to minimize the sum of squared errors that arises in the lot? That is, should you be a lumper or a splitter in your conception of the estimation problems you face? The surprise is that answering this question matters. In Section 11, we investigate whether this question has an objective answer.

### 8. The relation of shrinkage estimation to Bayesianism

So far our discussion has been in a frequentist framework. However, there is another prominent statistical framework, decision-theoretic Bayesianism, which conceives of estimation in a very different way. According to this framework, an agent should choose the estimate that has the lowest expected loss, where expected loss is calculated relative to the agent's posterior probability distribution.<sup>17</sup> Is shrinkage estimation and the holistic pragmatism it underwrites compatible with this Bayesian understanding of estimation? Our answer to this question is that it all depends on which type of Bayesianism is at issue.

The most minimal type of Bayesianism says that all it takes for an agent's probability distribution to be rational is synchronic coherence. A rational agent can have any probability distribution at any given time, so long as it obeys the axioms of probability. There is no further constraint on how prior probability values are assigned to different propositions, nor is there any diachronic constraint on how the agent's distributions at different times must be related. Minimal

<sup>&</sup>lt;sup>16</sup> For example, see the reply to Rudner (1953) by Jeffrey (1956).

<sup>&</sup>lt;sup>17</sup> More precisely, the expected loss of estimate *e* relative to distribution  $p(\theta | x)$  and loss function *L* is given by the formula  $\sum_{\theta} p(\theta | D) L(\theta, e)$ , where the sum is over all possible values of the parameter  $\theta$ .

Bayesians are therefore free to embrace shrinkage, since they can often engineer their probability distribution in such a way that a shrinkage estimator will end up having the lowest expected loss.

There are two stronger forms of Bayesianism that we think conflict with shrinkage estimation. Both embrace synchronic coherence but insist that rationality demands something more. The first interprets probabilities as rational degrees of belief and insists that prior probabilities should, in some sense, accurately and reasonably reflect agents' background knowledge or their "initial state" of information. The second type of Bayesianism updates probabilities by strict conditionalization or by Jeffrey (1983) conditionalization. Of course, these two types of Bayesianism are not mutually incompatible, and flesh and blood Bayesians often sign up under both banners.

To see the difficulties posed for Bayesians of the first kind, note that, if probabilities represent the degrees of belief of an agent, then an agent should never lump together estimation problems that the agent thinks are completely unrelated and independent of each other, even if the goal is to maximize global accuracy. This is because, for the kind of Bayesian we now are discussing, the estimate e' of the parameter  $\theta$ ' can affect the estimate e of a different parameter  $\theta$  only if there is some possible data x such that  $p(\theta | \theta' \& x) \neq p(\theta | x) - i.e.$ , only if it is possible that  $\theta$ ' could provide information about  $\theta$  over and above the information provided by the data itself.<sup>18</sup>

It may be helpful to demonstrate why this is the case, so we illustrate it for the case of just two parameters, although it holds in general. Suppose your goal is to maximize global accuracy over  $\theta$  and  $\theta$ ' in a Bayesian framework. Hence, you want to choose the estimates e of  $\theta$  and e' of  $\theta$ ', given data D, that jointly *minimize* the posterior expected total loss where this total loss is the sum of the loss on each of the parameters. In other words you want to pick the e and e' that jointly minimize:

$$\sum_{\theta,\theta'} p(\theta,\theta'|\mathsf{D}) * [\mathsf{L}(\theta,\mathsf{e}) + \mathsf{L}'(\theta',\mathsf{e}')]$$

<sup>&</sup>lt;sup>18</sup> Note that this is a purely synchronic constraint on the conditional probability distribution.

But if  $p(\theta | \theta' \& x) = p(\theta | x)$  for all x, then simple algebra and applications of the probability axioms show that the above formula reduces to:

$$\sum_{\theta} p(\theta|D)L(\theta, e) + \sum_{\theta'} p(\theta'|D)L'(\theta', e')$$

Since these two summands have no common variables, minimizing their sum just reduces to the problem of minimizing each of the summands. Hence, if  $\theta$  and  $\theta$ ' are thought by the Bayesian agent to be completely informationally independent in the sense that  $p(\theta | \theta' \& x) = p(\theta | x)$  for all x, then the problem of minimizing global accuracy just reduces to the problem of minimizing inaccuracy over each of the parameters, which means that there is no gain in lumping together estimation problems that are thought by the agent to be unrelated and independent, even if the goal is maximal global accuracy.

On the other hand, as we have argued in this paper, the holistic pragmatism underwritten by Stein's result sometimes recommends that agents who want to maximize global accuracy lump together unrelated and independent estimation problems. Hence, holistic pragmatism is clearly in conflict with a decision theoretic Bayesianism that interprets prior and posterior probabilities as representing the rational degrees of belief of some agent. Bayesians who want to embrace shrinkage estimation must consequently reject the view that the probabilities they use always represent their rational degrees of belief. This is not to deny that Bayesians can mathematically accommodate shrinkage estimators by, for example, using so-called "shrinkage priors," which are priors that are explicitly designed to induce Steinian shrinkage, as described in Efron (2011, pp 2-6), Efron and Morris (1973), or Lehmann (1983, p.299). But such "shrinkage priors" cannot reasonably be interpreted as representing the rational degrees of belief of some agent.

We have heard it said that Bayesians automatically incorporate shrinkage and that Bayesians therefore don't need to worry about artificially inducing shrinkage in order to maximize global accuracy. However, this is not correct. It is true that all proper priors – that is, priors that sum to 1 - impose a kind of shrinkage. In particular, any proper prior will "shrink" the estimate of a parameter towards values of the parameter that have a high prior probability. However, *Stein* shrinkage only happens when multiple parameters are shrunk towards a common point. The only way for a Bayesian to induce this kind of shrinkage is by imposing a common proper prior over all the parameters that are being estimated that has the effect of shrinking all the parameters towards the same point. If priors are interpreted as degrees of belief, however, then any such prior will represent a prior belief that the parameters are actually somewhat close together. Of course, sometimes such a prior belief is warranted, but our point here is that holistic pragmatism says that if your goal is to maximize global accuracy, you should *always* shrink the parameters towards a common point, regardless of whether you actually believe the parameters are close together. Now, as we saw in Section 4, the practical benefit of shrinkage will not be large unless the parameters are actually close to each other, but this does not alter the fact that the expected benefit of shrinkage is always greater than 0, provided that global accuracy is the goal.

Note, however, that if local accuracy rather than global accuracy is your goal, then Stein's result does not give Bayesian agents any reason for imposing a common prior over parameters that they think are independent. Thus, whether you should impose a common prior over the parameters you are estimating or keep the parameters separate depends on your goals. This observation brings us to the tension that we find with the second type of Bayesianism in which there is a commitment to updating by strict conditionalization or by Jeffrey conditionalization. The impediment here is that Stein's results show that you can and should change your estimates simply because your goals change. To see the problem, consider Ms. Multi-Tasker and Mr. One-at-a-Time as two stages in a single Jekyll-and-Hyde personality – namely you. On Monday you want to reduce your risk in estimating each of three parameters, while on Tuesday you want to reduce your total risk. MLE is admissible on Monday but not on Tuesday. Your estimates change value but not because of any new evidence you acquired. Updating by conditionalization is a way to take new evidence into account. It does not allow your estimates to change merely because you have changed your goals.

Of course, the Bayesian framework does allow you to change your estimate if your loss function changes, even if you do not gain any new evidence. But that is not what's going on in the above example. Even if your loss function remains the same throughout the whole week (e.g. you are using squared loss), Stein's result gives you a reason for using one estimate on Monday and another estimate on Tuesday. But if your loss function stays the same and you don't gain any new evidence, then the only way you can change your estimate in the Bayesian framework is by changing your probability distribution in a way that violates conditionalization. That is why shrinkage estimation and holistic pragmatism are in tension with the second type of Bayesianism.

We think the most plausible Bayesian response to Stein's results is to either reject them outright or to adopt an instrumentalist view of personal probabilities. The instrumentalist response is to abandon the idea that probabilities are rational degrees of belief. Rather, they are sometimes irrational (or arational), but agents should adopt them anyhow because they help agents get what they want. The former option, of outright rejection, has several motivations. The most obvious motive is perhaps that Stein's result is fundamentally a *frequentist* result since it is couched in terms of expected inaccuracy over all possible data sets. Bayesians might insist that Bayesians should not care about frequentist risk, but we know of few Bayesians who hold this view.<sup>19</sup>

Even if Bayesians grant that frequentist risk is relevant they still might insist that you shouldn't use frequentist risk, by itself, to decide which estimator to use. A more fully Bayesian decision theoretic solution would require that you calculate the expected utility of choosing each candidate estimator by averaging the frequentist risk of the various candidate estimators over a prior probability distribution. An estimator that maximizes expected utility relative to a particular prior is known as a "Bayes rule" relative to that prior. As it happens, MLE and shrinkage estimation both are Bayes rules relative to the (improper) prior that is flat over all the parameters you are estimating. In other words, relative to the flat prior, MLE and shrinkage estimation have the same expected utility. Consequently, if you come to the table with a flat prior in hand, you apparently have no reason (from a Bayesian point of view) for preferring shrinkage estimation to MLE.<sup>20</sup>

<sup>&</sup>lt;sup>19</sup> Angers and Berger (1985, p. 5) emphasize that frequentist risk should be important even to pure Bayesians because frequentist risk gives an indication of the average posterior expected loss.

<sup>&</sup>lt;sup>20</sup> We thank Teddy Seidenfeld for pressing us on this point. We note that this point only makes sense given certain presuppositions. In particular, countable additivity must be discarded.

We grant that an expected utility calculation will not tell you that shrinkage is preferable to MLE, if you use a flat prior in calculating expected utility. However, it's also true that such a calculation will not tell you that MLE is preferable to shrinkage, given that MLE and shrinkage have the same expected utility. If you want a non-arbitrary way of picking an estimator, you therefore need a tie-breaking criterion aside from expected utility. A reasonable tie-breaker, we think, is to pick the estimator that has lower frequentist risk.

The above discussion assumes that you have already adopted a flat prior. But what if you are trying to decide what prior to adopt in the first place? An expected utility calculation cannot tell you what prior to adopt since expected utility must be calculated relative to a prior. If you are unsure of what prior to adopt, but your goal is to minimize global error, then Stein's result gives you a reason for preferring a prior that imposes shrinkage to a prior that does not.

Some Bayesians argue that shrinkage estimators should not be used in the absence of genuine prior information and in particular that "shrinkage priors" of the sort described in this section should never be used. For example, Angers and Berger (1985) show that under certain (strong) assumptions, Bayesians who are certain or nearly certain that several estimation problems are probabilistically independent should not combine the problems and impose shrinkage because it is possible to obtain more robust estimates (i.e., estimates that are less sensitive to the choice of prior) by considering the problems separately.<sup>21</sup> We concede that shrinkage priors do not provide you with robust individual estimates under the assumptions that Angers and Berger describe. Indeed, as we discuss further in the next section, what point you shrink towards (or what shrinkage prior you choose to adopt) can greatly affect which of your individual estimates end up being more accurate and which ones end up being less accurate when compared to estimates provided by MLE. Nonetheless, if your main concern is to maximize global accuracy – as opposed to making sure that you have robust estimates – then Stein's result shows that you ought to shrink.

<sup>&</sup>lt;sup>21</sup> Perlman and Chaudhuri (2012) offer a different argument for a similar conclusion. They claim, without offering any explanation, that agents who use shrinkage estimation in the absence of prior information will unwittingly end up using a procedure that has the effect of reversing the Stein effect (see n24 for a description of the procedure).

There are two remaining and perhaps more fundamental reasons why Bayesians might reject the use of shrinkage estimation. One reason is that shrinkage estimation is a type of *estimation*, and some Bayesians have claimed that estimation is simply not that interesting from a Bayesian perspective. Another fundamental reason that Bayesians might want to reject the use of shrinkage estimation is that they value language invariance, and shrinkage estimators are not language invariant; we discussed language invariance in Section 3.

### 9. Shrinkage Estimation, AIC, and the Bias/Variance Tradeoff

Complex models (ones with more adjustable parameters) will in general exhibit less *estimation error* or *bias* than simpler models (ones with fewer parameters) because the extra flexibility that accompanies complexity enables complex models to accurately fit real patterns in the data.<sup>22</sup> For example, an  $n^{\text{th}}$  degree polynomial can fit a data set containing n-1 observations perfectly whereas polynomials with fewer adjustable parameters will almost never do as well. Unfortunately, complexity comes at a cost; more complex models are also more likely than simpler ones to be misled by noise. Hence, more complex models have greater *approximation error* or *variance* since a complex model is likely to "bounce around" quite a bit when it is fitted to different data sets drawn from the same underlying distribution.

The Akaike Information Criterion (AIC) gives advice concerning how bias should be traded off against variance when the goal is to maximize predictive accuracy. According to AIC, the predictive inaccuracy of model M given data D can be estimated by the number of parameters in M minus the log-likelihood of the best-fitting member of M. AIC tells you that the model that has the lowest AIC score will be the one that is most predictively accurate; this is the

<sup>&</sup>lt;sup>22</sup> In speaking of the "bias" and "variance" of models and estimators, we are following the (perhaps unfortunate) statistical practice of using these terms ambiguously. In the context of parameter estimation, "bias" and "variance" have precise technical meanings, as we also note in the text. In model selection, on the other hand, "bias" just means something like "the inability of a model to mimic the true curve, whatever the true curve happens to be." For example, LIN is "biased" in this latter sense because it can adequately mimic the true curve only if the true curve happens to be roughly linear.

model that achieves the best balance between fit and complexity, or, in other words, between bias and variance. The reason scientists often should prefer simple models that they know are false (i.e., ones that are biased away from the truth) over more complex models that they know to be true is that simpler models have lower variance and hence often have a higher degree of predictive accuracy (see Forster and Sober 1994).

Shrinkage estimators rely on a similar kind of bias-variance trade-off. However, in the context of parameter estimation, the terms "bias" and "variance" take on technical meanings that differ from the meanings of those terms in model selection. For example, an unbiased estimator is one whose expected value equals the true value of the estimated parameter. Using these technical concepts makes it possible to decompose the error of an estimator in a precise way. In particular, if you look at a single estimator, m, of some quantity  $\theta$ , you can decompose its mean squared error in the following way (see e.g., Wasserman 2004, p. 91):

$$MSE(m) = [bias(m)]^2 + variance(m).$$

If m is an ML estimator, then  $[bias(m)]^2 = 0$  when the distribution of the data is normal, since MLE is unbiased in this case. Hence, in this case, the MSE of m simply reduces to the variance of m. If we *independently* estimate several parameters  $\theta_1, \theta_2, ..., \theta_n$  and obtain  $m_1, m_2, ..., m_n$  as their ML estimates, then the total MSE of all the estimates is simply the sum of all their variances:

 $MSE(m_1, m_2, ..., m_n) = variance(m_1) + variance(m_2) + ... + variance(m_n).$ 

Since using a shrinkage estimator yields a lower total MSE than MLE, and since shrinkage estimators are clearly biased (as noted, the ML estimator is unbiased, and shrinkage means shrinking your estimates *away* from the ML estimates; note that we are here using the terms "biased" and "unbiased" in the technical sense), it is clear that shrinkage estimators work because they lead to a reduction in overall variance. By shrinking all the ML estimates toward some common point, you bias your estimate of each parameter, but at the same time you reduce the freedom of each estimate to vary in response to noise; by "tethering" all the estimates to a

single point, shrinkage thereby reduces the overall variance in your estimates. Just as AIC sometimes prefers simpler false models because they have lower variance, so shrinkage estimators sacrifice the unbiasedness of MLE for the sake of obtaining a lower variance.

Trading-off bias against variance arises in many statistical contexts<sup>23</sup> – in model selection, in classification problems (see, e.g. von Luxburg and Scholkopf 2009, pp. 662-664) and – as we have seen – in parameter estimation. What is especially striking about shrinkage estimators is that there are several conflicting ways of sacrificing bias to achieve a reduction in variance; recall that the James-Stein estimator shrinks all estimates towards zero while the Efron-Morris estimator shrinks them towards the grand sample mean. However, whether you introduce *this* bias or *that* bias is less important than the fact that you introduce *some bias or other*.<sup>24</sup> What bias you introduce (i.e. what the point is towards which you shrink) will have consequences for which individual estimates end up being less accurate than they might have been had you used MLE, but *overall* the reduction in variance justifies sacrificing accuracy at the level of estimating individual parameters.

# 10. Shrinkage, Realism, and Instrumentalism

Model selection criteria like AIC legitimate a kind of *instrumentalism* (Sober 2008). According to this interpretation, the goal of AIC is not to decide which model is true, or has the highest probability of being true, but rather to determine which model (among the candidate models considered) will make the most accurate predictions. As noted, a surprising property of

<sup>&</sup>lt;sup>23</sup> An anonymous reviewer pointed out to us that similar trade-offs also arguably occur outside of statistics, e.g. in the "Runge phenomenon" in numerical analysis.

<sup>&</sup>lt;sup>24</sup> It is worth noting that *how* you choose to bias your estimates is important. Perlman and Chaudhuri (2012) show that there are procedures for picking the point towards which you shrink that lead to a "reverse" Stein effect wherein the resulting shrinkage estimator does worse than MLE in expectation. In particular, the point you shrink towards needs to be relatively stable given different data sets; otherwise, your shrinkage estimator is not going to reduce total variance. Here's an example: given data  $(x_1, x_2, ..., x_n)$  about parameters  $(\theta_1, \theta_2, ..., \theta_n)$ , consider  $(x_1, x_2, ..., x_n)$  as the center of a sphere in n-dimensional space and then randomly pick some point within the sphere towards which you shrink your data. This procedure gives you an estimator that bounces around given different data sets, and that therefore doesn't help you reduce overall variance.

AIC is that it sometimes (correctly) judges that a model known to be false will be more predictively accurate than a model known to be true. Should shrinkage estimators be placed under the same instrumentalist umbrella? According to this interpretation, the goal is not to discover the true value of a parameter, but to find the estimate that will most accurately predict new data. The surprise is that achieving this goal depends on whether you want to predict new observations of one or two parameters, or of three or more. And if you want to predict three or more, a second surprise presents itself; it matters whether you want to reduce total expected error or rather want to minimize expected error on each parameter.

Although this instrumentalist gloss makes sense, shrinkage estimation is something that realists also need to take on board.<sup>25</sup> The realist view of estimation is that the goal is to get as close as possible to the true value of the quantity being estimated. Shrinkage estimators achieve that goal better than straight MLE when there are more than three parameters being estimated and the goal is to minimize total inaccuracy. However, some realists will be lumpers while others will be splitters, and so they will disagree about how Stein's result should bear on their scientific practice.

### 11. When to lump and when to split?

Can the question of when to lump and when to split be settled in an objective way? Lumping the success rates of baseball players is a good idea if the goal is to minimize *global* inaccuracy. However, as we emphasized earlier, Stein's result does *not* guarantee that shrinkage estimation will do better than MLE when you treat each baseball player as a separate estimation problem. As noted, MLE is admissible when you estimate the value of a single parameter. This does not mean that MLE dominates shrinkage estimators when k<3; it means only that no shrinkage estimator dominates MLE. Furthermore, the fact that no shrinkage estimator dominates MLE in single-parameter estimation just means that there is no shrinkage estimator that does at least as well as MLE across *all* possible values that the parameter might have; this

 $<sup>^{25}</sup>$  The same point holds for AIC – realists can embrace this estimator so long as closeness to the truth is understood in the right way (Sober 2015).

leaves open that there may be estimators that do significantly better than MLE when the parameter has some *particular* true value. Indeed, if you return to Figure 1, you'll notice that dividing the sample mean in half actually has lower risk than the sample mean itself when the true value of the parameter is sufficiently close to 0. "Admissible" does not mean "optimal in all cases."

We now change our focus to a question that the concept of admissibility cannot answer: should you be a lumper or a splitter in the way you formulate your estimation problems? In general, there are two things that can happen when you lump three or more estimation problems and use shrinkage, rather than split them and use MLE:

- (1) The risk of some of the individual estimates increases substantially.
- (2) The risk of each of each individual estimate remains about the same or decreases.<sup>26</sup>

If (1) holds, you face a dilemma; shrinkage estimation will decrease your global risk, but only at the price of increasing the risk of some of the individual estimates. Should you trade an increase in risk at the individual level for a reduction in risk at the global level? The answer to this question depends on what your goals are and what kind of error you prefer to avoid.<sup>27</sup>

On the other hand, if (2) holds, you objectively ought to use a shrinkage estimator; by shrinking, you'll reduce global risk without increasing individual risk, or you'll reduce both. In

<sup>&</sup>lt;sup>26</sup> It may seem that (2) contradicts the fact that MLE is admissible when you are estimating a single parameter. However, the fact that MLE is admissible just means that, given (normally distributed) data D and a single parameter p, there is no function of D that weakly dominates the ML estimator. This does not exclude the possibility that you can get a better estimate of p by lumping D with another data set D' and using a shrinkage estimator on D&D'.

<sup>&</sup>lt;sup>27</sup> The measure of global inaccuracy that we have relied on so far in this paper implicitly places an equal weight on each of the individual estimation problems, since global inaccuracy is simply the *unweighted* sum of the inaccuracies in all the individual estimates. It may, however, happen that getting accurate estimates for some of the parameters is more important than getting accurate estimates for others. Brown (1975) models this by weighting each inaccuracy term  $(x_i-\theta_i)^2$  by a factor  $c_i$  that measures the importance of getting an accurate estimate of  $\theta_i$ , and he shows that Stein's result is surprisingly robust across different assignments of weights (although he does not propose explicit estimators corresponding to the different weightings of the local estimation problems). In a similar vein, Efron and Morris (1972) introduce "compromise estimators" that aim at lowering total inaccuracy while at the same time limiting the loss in accuracy at the level of individual estimates.

general, (2) holds if and only if the true means and variances of the different parameters you are estimating are sufficiently close to each other. To see why, recall that shrinkage estimators work by introducing a bias in order to reduce variance. If the true means are sufficiently close to each other, then the bias introduced by using a shrinkage estimator such as the Efron-Morris estimator will be smaller than the reduction in variance *even at the level of individual estimates*.

How do you know whether the quantities you are studying have true means and variances that are "sufficiently close" before you have the relevant data? Sometimes background knowledge can serve as a guide. Suppose, for example, that you wish to estimate the mean heights in all the northern European countries. Your background knowledge might lead you to believe that the true mean heights and variances for the different countries, whatever they are, are close together. If your belief is correct, each true population mean will be close to the true grand mean, which entails that the bias introduced by the Efron-Morris estimator will be small for each of the individual estimates you make. At the same time, combining your samples from several countries in the way that the Efron-Morris estimator does will reduce the total variance of your estimates, as we discussed in Section 9. Lumping therefore seems justified in this case; you have reason to believe that (2) is true.

Background knowledge would also seem to justify lumping together Major League baseball players. Although the players do vary in ability, there is reason to believe that their true abilities are sufficiently close so that by lumping the estimation problems together, you can expect to increase the accuracy in your estimate of *almost every* player. Of course, your estimate of players who are truly outstanding (i.e. whose true batting success probability is far from the grand mean) will suffer, but if you know who those players are, you can just exclude them from shrinkage.

These examples show how using background knowledge to lump together estimation problems can be expected to produce better individual estimates if they are done right, but the examples offer no general guidelines for how to do this. According to Efron, "[s]cientific guidance would be most welcome at this point…" (2011 p. 185). Efron goes on to quote Miller's (1981) pronouncement (in a slightly different context) that this is where "statistics takes leave of

mathematics and must be guided by subjective judgment." Efron explains that subjective judgment is still what dominates statistical practice, but goes on to offer "hints ... of a more principled approach to separation and combination ... but at this stage they remain just hints" (p. 186).

The upshot is that lumping together several estimation problems and using a shrinkage estimator is sometimes *objectively better* than keeping the problems separate and using MLE, in the sense that lumping can sometimes yield individual estimates that are more accurate than the estimates you would get by handling the problems separately. This situation arises when the populations under study are "similar enough." But when and how problems should be lumped together or split apart remains an important open problem in statistics.<sup>28</sup>

## **12. Concluding Comment**

Stein's result seems paradoxical when estimating a quantity is taken to be closely connected to assigning a probability to a proposition.<sup>29</sup> When propositions A, B, and C are probabilistically independent of each other, the probability of their conjunction is determined by the probabilities of the conjuncts:

<sup>&</sup>lt;sup>28</sup> There is a parallel situation for AIC. Consider NULL and DIFF as claims about how the mean heights in two (or more) human populations are related. NULL says they have the same mean height; DIFF says that the heights may differ. Since NULL says there is a single mean height that characterizes each population, it has a single adjustable parameter. DIFF has one adjustable parameter for each population. AIC tells you to prefer NULL if the sample means are close together and DIFF when they are very far apart, where the question "how close is close enough?" is answered by considering how the two models differ in their numbers of adjustable parameters. NULL lumps whereas DIFF splits.

<sup>&</sup>lt;sup>29</sup> This is not to say that the puzzling quality of Stein's result derives entirely from the assumption that estimating a quantity and assigning a probability to a proposition are closely connected projects.

If Pr(A)=x and Pr(B)=y, Pr(C)=z, and A,B,C are probabilistically independent of each other, then Pr(A&B&C) = xyz.

This is true by definition. Analogy might suggest that something similar is true of estimation. One might expect, when the quantities Q, R, and S are probabilistically independent of each other, that the following principle holds:

If BE(Q)=x, BE(R)=y, BE(S)=z and Q,R,S are probabilistically independent of each other, then BE<Q,R,S>=<x,y,z>.

Here BE( $\theta$ ) is the best estimate of  $\theta$ , given the data at hand, and "<...>" denotes a vector of variables or values. This principle says that the best estimate for a vector of quantities is determined by the best estimate for each quantity when the quantities are probabilistically independent. The principle would be right if the best estimate of a quantity were the one that has the highest probability (density) of being true. This is an intuitive interpretation of estimation, but it is not the one that frequentists embrace. Their interpretation is built on the concept of admissibility. In this paper, we have attempted to describe the surprising form of holistic pragmatism that issues from that conception.<sup>30, 31</sup>

#### References

Angers, Jean-Francois and Berger, James O. (1985) "The Stein Effect and Bayesian Analysis: A Reexamination," *Technical Report #85-6*. Department of Statistics, Purdue University.

<sup>&</sup>lt;sup>30</sup> In thinking about how assigning a probability to a proposition is related to estimating a quantity, it is important to note that an estimation problem must involve infinitely many possible values if shrinkage estimators are to have lower expected error than MLE (Guttmann 1982).
<sup>31</sup> We thank Marty Barrett, Larry Brown, Jan-Willem Romeijn, Teddy Seidenfeld, Mike Steel, Reuben Stern, and the anonymous referees for very useful comments. This paper is dedicated to the memory of Charles Stein (1920-2016).

Baranchik, Alvin J. (1964) "Multiple Regression and Estimation of the Mean of a Multivariate Normal Distribution." *Technical Report 51*. Department of Statistics, Stanford University.

Blyth, Colin (1951) "On Minimax Statistical Decision Procedures and their Admissibility." *The Annals of Mathematical Statistics*, Vol. 22 (1): 22-42.

Bock, Mary E. (1975) "Minimax Estimators of the Mean of a Multivariate Distribution." *Annals of Statistics*, Vol. 3 (1): 209-218.

Brown, Lawrence D. (1966) "On the Admissibility of Invariant Estimators of One or More Location Parameters." *Annals of Mathematical Statistics*, Vol. 37 (5): 1087-1136.

Brown, Lawrence D. (1971) "Admissible Estimators, Recurrent Diffusions, and Insoluble Boundary Value Problems." *The Annals of Mathematical Statistics*, Vol. 42 (3): 855-903.

Brown, Lawrence D. (1975) "Estimation with Incompletely Specified Loss Functions (the Case of Several Location Parameters)." *Journal of the American Statistical Association*, Vol. 70 (350): 417-427.

Carnap, Rudolf (1950) "Empiricism, Semantics, and Ontology." *Revue Internationale de Philosophie*, Vol. 4 (2): 20-40.

Edwards, A. W. F (1974) "The History of Likelihood." *International Statistical Review*, Vol. 42 (1): 9-15.

Efron, Bradley (2013) Large-Scale Inference: Empirical Bayes Methods for Estimation, Testing, and Prediction. Cambridge University Press.

Efron, Bradley and Morris, Carl (1972) "Limiting the Risk of Bayes and Empirical Bayes Estimators – Part II: The Empirical Bayes Case." *Journal of the American Statistical Association*, Vol. 67 (337): 130-139. Efron, Bradley and Morris, Carl (1973) "Stein's Estimation Rule and Its Competitors – An Empirical Bayes Approach." Journal of the American Statistical Association, Vol. 68 (341): 117-130.

Efron, Bradley and Morris, Carl (1977) "Stein's Paradox in Statistics." *Scientific American* Vol. 236 (5): 119–127.

Forster, Malcolm and Sober, Elliott (1994) "How to Tell When Simpler, More Unified, or Less Ad Hoc Theories Will Provide More Accurate Predictions." *British Journal for the Philosophy of Science*, Vol. 45: 1-36.

Galton, Francis (1888) "Co-relations and their Measurement, chiefly from Anthropometric Data." *Proceedings of the Royal Society of London*, Vol 45: 135-45.

Gauss, Carl F. (1823) "Theoria Combinationis Observationum Erroribus Minimis Obnoxiae: Pars Posterior." Translated (1995) as *Theory of the Combination of Observations Least Subject to Error: Part One, Part Two, Supplement*. Translated by: G. W. Stewart. Society for Industrial and Applied Mathematics.

Gruber, Marvin (1998) Improving Efficiency by Shrinkage: The James-Stein and Ridge Regression Estimators. CRC Press.

Guttmann, Sam (1982) "Stein's Paradox is Impossible in Problems with Finite Sample Space." Annals of Statistics, 10(3): 1017-1020.

Hodges, Joseph and Lehmann, Erich (1951) "Some applications of the Cramér-Rao inequality."*Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability*,Berkeley and Los Angeles, University of California Press, pp. 13-22.

James, Willard and Stein, Charles (1961) "Estimation with Quadratic Loss." *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability*, Vol. 1: 361-379.

James, William, (1896, 1979) "The Will to Believe" in F. Burkhardt et al. (eds.), *The Will to Believe and Other Essays in Popular Philosophy*, Cambridge: MA, Harvard, pp. 291–341.

Jeffrey, Richard (1956) "Valuation and Acceptance of Scientific Hypotheses." *Philosophy of Science*, Vol. 23 (3): 237-246.

Jeffrey, Richard (1983) *The Logic of Decision*. Second edition. Cambridge University Press, Cambridge.

von Luxburg, Ulrike and Schölkopf, Bernhard (2009) "Statistical Learning Theory: Models, Concepts, and Results." In: D. Gabbay, S. Hartmann, and J. Woods (Eds). *Handbook of the History of Logic*, Vol 10: Inductive Logic.

Miller, Rupert G. Jr. (1981) Simultaneous Statistical Inference: Second Edition. Springer.

Pascal, Blaise (1662) *Pensées*. Translated by W. Trotter, New York: J. M. Dent Co., 1958, fragments: 233-241.

Perlman, Michael D. and Chaudhuri, Sanjay (2012) "Reversing the Stein Effect." *Statistical Science*, Vol. 27 (1): 135-143

Rudner, Richard (1953) "The Scientist *Qua* Scientist Makes Value Judgments." *Philosophy of Science*, Vol. 20 (1): 1-6.

Sober, Elliott (2008) *Evidence and Evolution – the Logic Behind the Science*. Cambridge University Press.

Sober, E. (2015) *Ockham's Razors – A User's Manual*. Cambridge: Cambridge University Press.

Spanos, Aris (2016) "How the Decision Theoretic Perspective Misrepresents Frequentist Inference: `Nuts and Bolts' vs Learning from Data" Available at https://arxiv.org/pdf/1211.0638v3.pdf

Stein, Charles (1956) "Inadmissibility of the Usual Estimator for the Mean of a Multivariate Normal Distribution." *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability*, Vol. 1: 197-206.

Stein, Charles (1962) "Confidence Sets for the Mean of a Multivariate Normal Distribution (with discussion). *Journal of the Royal Statistical Society. Series B (Methodological)*, Vol. 24 (2): 265-296.

Stigler, Stephen (1990) "The 1988 Neyman Memorial Lecture: A Galtonian Perspective on Shrinkage Estimators." *Statistical Science*, Vol 5 (1): 147-155.

Strawderman, William E. (1971) "Proper Bayes Minimax Estimators of the Multivariate Normal Mean." *Annals of Mathematical Statistics*, Vol 42 (1): 385-388.

Wasserman, Larry (2004) All of Statistics: A Concise Course in Statistical Inference. Springer.

White, Morton (2005) A Philosophy of Culture: The Scope of Holistic Pragmatism. Princeton University Press.