

There is Individualized Treatment. Why Not Individualized Inference?

Keli Liu and Xiao-Li Meng
(*keliliu@stanford.edu* & *meng@stat.harvard.edu*)
Stanford University and Harvard University

November 29, 2015

Abstract

Doctors use statistics to advance medical knowledge; we use a medical analogy to introduce statistical inference “from scratch” and to highlight an improvement. Your doctor, perhaps implicitly, predicts the effectiveness of a treatment for *you* based on its performance in a clinical trial; the trial patients serve as *controls* for you. The same logic underpins statistical inference: to identify the best statistical procedure to use for a problem, we simulate a set of *control problems* and evaluate candidate procedures on the controls. Now for the improvement: recent interest in personalized/individualized medicine stems from the recognition that some clinical trial patients are better controls for you than others. Therefore, treatment decisions for you should depend only on a subset of *relevant* patients. *Individualized statistical inference* implements this idea for control problems (rather than patients). Its potential for improving data analysis matches personalized medicine’s for improving healthcare. The central issue—for both individualized medicine and individualized inference—is how to make the right *relevance robustness trade-off*: if we exercise too much judgement in determining which controls are relevant, our inferences will not be robust. How much is too much? We argue that the unknown answer is the Holy Grail of statistical inference.

Prologue: The Data Doctor

A usual course on statistical inference teaches its *mechanics*: how to construct confidence intervals or conduct hypothesis tests through the use of probabilistic calculations. We then become so fluent in the modern language of statistics that we forget to ask: Why did the language develop this way? Can inference be done using a different language? For example, why introduce the notion of probability at all—is it truly indispensable? This article does not aim to introduce the reader to the language of inference as currently spoken by many statisticians. Instead, we try to create a language for inference “from scratch,” and stumble upon the main statistical dialects (and a few variants) by accident. By from scratch, we mean that we justify each step in the construction employing only “common sense”. In fact, we claim that one only needs to understand the following problem to understand statistical inference.

The Doctor’s Problem. A doctor needs to choose a treatment for her patient, Mr. Payne. Her options are the standard treatment, *A*, and an experimental treatment, *B*. What does the doctor do to make her decision? She goes and finds out how *A* and *B* worked on *other* patients. Suppose she discovers that treatment *B* outperformed *A* on patients in a large randomized clinical trial. She wants to apply this result to Mr. Payne. For this to be a reasonable action, the patients in the clinical trial need to be good *controls* for Mr. Payne (who is 50 years old, weighs 200 pounds, exercises once a month, etc.). Of course, she realizes that not all patients in the trial are good controls (certainly not the 12 year old girl). So she selects a subset of patients

who closely match Mr. Payne’s characteristics and bases her decision on A and B ’s performance over this subset.

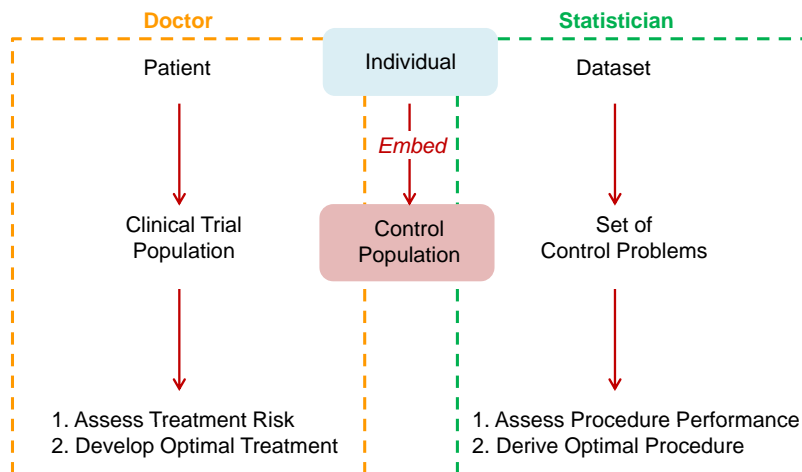


Figure 1: The statistician needs to treat his “patient”, a dataset, with a statistical procedure. He chooses between candidate procedures based on their performance across a set control problems.

In summary, the doctor does two things: (i) she obtains controls for Mr. Payne and (ii) among the initial controls, she selects the subset that is most relevant for Mr. Payne. Statisticians do the same two things; we are the doctors for data (see Figure 1). Throughout this review, we exploit this analogy to derive insights into statistical inference. A statistician’s patients are datasets for which he hopes to choose the analysis procedure leading to the most accurate scientific conclusions. To test the effectiveness of candidate procedures, a statistician must also set up a clinical trial; he needs to simulate *control problems* and see how well his procedures perform on the controls. A great many jobs involve doing something similar: a rocket scientist runs simulations before launching an actual rocket, a basketball coach runs scrimmages to simulate a real game. But just as we wouldn’t trust a statistician to design rocket simulations, so there are tricks for getting simulated control problems to closely match the statistical problem of interest. That is, creating “good” control problems represents the key stumbling block in inference.

In our attempt to create good controls, we will encounter a fundamental tradeoff: *the relevance–robustness tradeoff*. It is easy to understand in the context of the doctor’s problem. If we estimate Mr. Payne’s treatment response by the average response across all clinical trial patients, we will likely do poorly; patients whose background variables differ from Mr. Payne’s may respond to treatment differently. But if we predict Mr. Payne’s treatment response using only those patients who match Mr. Payne’s background exactly, we will still do poorly. There are too few (or none) of such patients in the clinical trial; accordingly, our estimate will be too variable. Thus, our inferences can be hurt either by too much matching (not robust) or too little matching (not relevant). We will see that this lesson generalizes directly to the statistical context where we must decide how closely the features of a control dataset should match those of the actual dataset in order for the control problem to be deemed relevant.

No consensus exists on the “right” degree of matching, neither in the doctor’s case nor the statistician’s. However, it seems intuitively clear that no single answer can work for all problems. This maxim is the inspiration for personalized medicine: for each patient, the doctor begins *anew* the process of selecting controls, striking a balance between relevance and robustness that is tailored to the patient’s situation. Similarly, we should try to individualize statistical inferences to the problem at hand. We will argue that this spirit of individualization is counter to the two

standard statistical methodologies—Bayesian and (unconditional) Frequentist inference—which coincide with the *static* positions of always performing complete matching or never performing any matching. The need for personalized medicine is clear. We hope this article will make a convincing argument that the need for *individualized inference* is equally urgent.

The layout for the remainder of our article is as follows:

Section 1. A doctor can go out into the world and find controls for Mr. Payne. How can a statistician create control problems when he only has one data set? Anyone who has taken standardized tests knows that “problems” fall into types. If we can use the observed data to identify the “type” of a statistical problem, we can then simulate control problems of the *same type*.

Section 2. Simulated control problems will not perfectly resemble our actual problem. We should draw inferences using only those controls which match the actual problem on important features, e.g., the sample size. Through a series of examples, we show how to weigh gains in our inference’s relevance (from matching on a feature) against potential losses in its robustness.

Section 3. Many judgments underpin the selection of a set of “relevant” controls. We discuss strategies to make our inferences robust to *mis*judgments. These include the use of pivots to create confidence intervals and the minimax principle. We give a unifying view of these strategies as ways to “create symmetry/invariance.”

Section 4. We map different statistical methodologies according to the relevance-robustness tradeoff they adopt in creating control problems. We see that Bayesian inference coincides with complete matching and (unconditional) Frequentist inference with no matching. Neither extreme seems appealing; we suggest compromises based on partial matching, which include the controversial Fiducial inference.

1 What Makes It a *Statistical* Inference?

1.1 It’s the Built-in Uncertainty Assessment ...

We go to the doctor because, unlike Google, the doctor can distinguish (hopefully) which treatment among available options is most appropriate *for me*. Similarly, the statistician must choose the method of analysis most appropriate for a *particular* dataset. A common logic governs the doctor’s and statistician’s choice: to assess the effect of a procedure (treatment B) on a target individual (Mr. Payne), we need to obtain “controls” for that target individual.

Control Patients. Let \ominus denote Mr. Payne, y his health outcome under treatment B and \bar{x} his vector of observed background variables (age, weight, height, etc.). If we could clone Mr. Payne and apply treatment B to his clone, we would know y perfectly (a *deterministic* inference). Instead, we settle for a clinical trial of control patients, for whom we obtain data $\{(\bar{x}', y')\}$, where (\bar{x}', y') are the background variables and response for \ominus' . Throughout, we will use prime notation to denote aspects of controls and the unprimed versions of the same quantities to refer to aspects of the target individual. The prime notation reflects our desire for controls to mimic our target, \ominus , as closely as possible.

Control Problems. The target “individual” \ominus , is now a statistical “problem” which comprises two parts: (i) the dataset, D , being analyzed and (ii) some unknown “truth”, θ , we hope to infer, which may or may not correspond to a parameter in a statistical model. For concreteness, suppose we observe a sequence of 9 colored pixels (D) and wish to predict the color of the 10th pixel (θ), as in Figure 2b. There are a wealth of prediction algorithms we could use to accomplish this task. To choose among them, we evaluate their accuracies on a set of control problems (D', θ'). Each control problem is a sequence of 10 pixels. We apply the candidate prediction algorithms to the first 9 pixels (D'), and check their outputs against the 10th (θ').

In our analogy, the statistical problem, $\ominus = (D, \theta)$, plays the role of Mr. Payne, $\ominus = (\bar{x}, y)$. The doctor observes \bar{x} and wishes to infer/predict y ; the statistician observes D and wishes to infer θ . An important consequence of this parallelism is that just as we improve the relevance

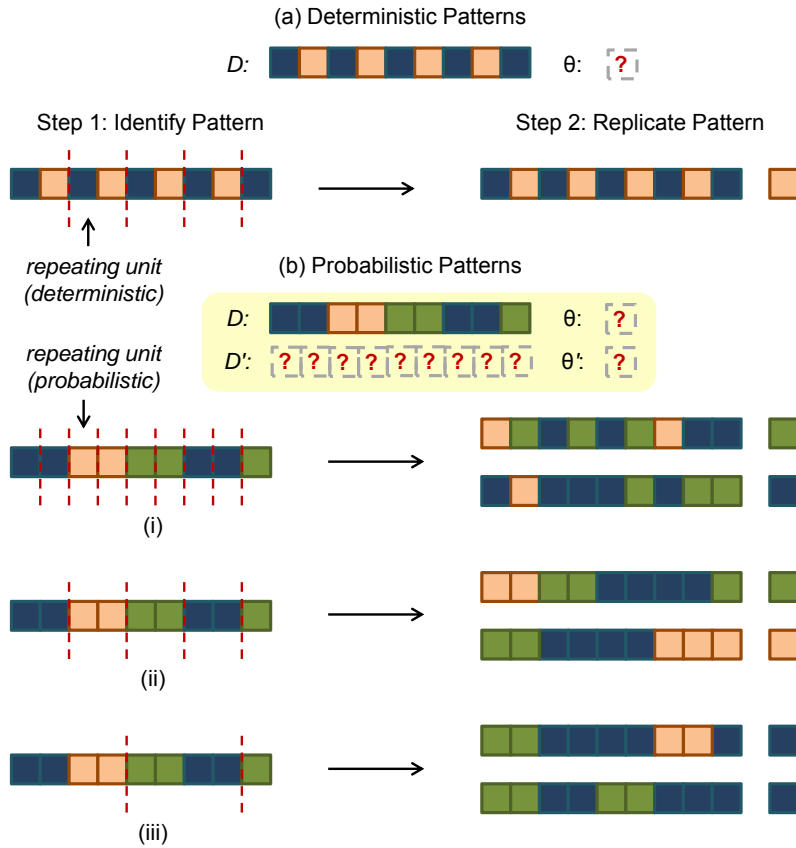


Figure 2: Diagram (a) shows a sequence following a *deterministic* pattern, and diagram (b) shows a sequence constructed from replications of a *probabilistic* pattern. To create controls, we first decompose the data into the repeating subunits of a pattern. Plots (i), (ii), and (iii) represent three possible decompositions. We then simulate controls by randomly sampling (with replacement) from the set of observed subunits and stringing the sampled subunits together until we obtain 10 pixels or more. To the right of each arrow are two example controls simulated according to the decomposition on the left.

of our inference for Mr. Payne by choosing controls which match his background, \vec{x} , we should ensure control problems, (D', θ') , for $\odot = (D, \theta)$ match its “background”, D (see Section 2.3).

A Conundrum: Creating Controls from D Alone. A doctor can go out into the world and solicit clinical trial subjects. Her controls therefore contain information *external* to Mr. Payne’s observed characteristics, \vec{x} . But typically the *only* input the statistician has in creating control problems is the data, D , supplemented by common sense and subject matter knowledge. It is absurd to think that one can evaluate the risks of various treatments based on observing Mr. Payne’s background, \vec{x} , alone, i.e. without experimentation on others. Since (D, θ) parallels (\vec{x}, y) in our analogy, isn’t it equally absurd to try and conduct inference using only D ? Up to this point, our reasoning process has closely paralleled the doctor’s. It is here that statistics makes a wholly original contribution in supplying a logic for creating meaningful controls (D', θ') from D alone. Not only do we possess the ability to predict θ from D , but we can also assess the accuracy of our prediction using only D . The capacity to assess uncertainty is *built into* the data. This is the *magic* of statistics. Of course, all magics come with “hidden designs,” which we will reveal below.

1.2 The Data As A Replicating Pattern

The pixel example (Figure 2) hints at a strategy for creating controls (D', θ') based only on D . To make it explicit, consider a simpler problem (depicted in Figure 2a) where we *assume* the sequence follows a deterministic pattern. From the observed data, we easily identify the repeating unit of this pattern as a block of blue-orange pixels. Using this fact, we can simulate controls which follow the same pattern as the target. At a high level, our strategy for creating controls can be summarized as follows:

1. Assume the target problem, (D, θ) , follows a pattern.
2. Extract/estimate the pattern using only the observed data, D .
3. Produce controls, (D', θ') , based on the extracted pattern.

For the deterministic pattern in Figure 2a, this strategy would lead to perfect controls: $(D', \theta') = (D, \theta)$.

Probabilistic Patterns. A *deterministic* pattern would not be useful in describing the (seemingly) more complex sequence in Figure 2b. So we extend our strategy: use a *probabilistic* pattern to describe the color scheme of the first 9 pixels and assume that the 10th pixel follows the same pattern. A probabilistic pattern is typically captured by a (probability) *distribution*—it is a pattern for *how things vary*. For example, when we say that height in a human population follows a normal distribution, we mean that *variation* in height follows the pattern of a bell shaped curve. In the current problem, we can try to describe the *variation* in pixel colors using a distribution, call it f . For example, based on the 9 observed pixels (4 blue, 3 green, 2 orange), we might guess that each individual pixel has probability $4/9$, $3/9$ and $2/9$ respectively of being blue, green or orange (independent of the other pixels). This constitutes an estimate \hat{f} of f . We can then produce control sequences, (D', θ') , according to the estimated pattern \hat{f} ; the two example controls on the right of Figure 2b(i) were simulated in this way.

Robustness versus Descriptive Power. The above pattern assumes that the colors of individual pixels are independent of each other. This pattern fails to capture the block structure of the sequence: the first two colors are the same, so are the next two, and so on. To fix this deficiency, we can choose *blocks* of pixels as our repeating unit, as in Figure 2b(ii) and (iii). That is, color variation, not for a single pixel but for a block of pixels, follows some pattern f . Our estimated pattern \hat{f} will then encode features of the observed sequence that are specific to pixels *as well as* those specific to blocks; see the controls corresponding to decompositions (ii) and (iii) in Figure 2b.

With this greater descriptive power comes the danger that our pattern captures features of the observed pixels that do not generalize (artifacts). Under decomposition (iii), we observe two blocks of size 4: a block of blue-blue-orange-orange (bboo) pixels and a block of green-green-blue-blue (ggbb) pixels. We might guess that each non-overlapping block of 4 pixels is bboo or ggbb with probability $1/2$. This choice of \hat{f} produces controls with a rather curious feature:

a two pixel orange block *must be* preceded by a two pixel blue block. Prediction algorithms which exploit this feature would do well on our simulated controls; they would do poorly on our actual problem if such a feature turns out to be an artifact. The point is, we have no way to assess whether such a feature is genuine or not since a block of orange pixels appears only once in the observed data. If we had further replications of 4 pixel-blocks in our observed data, we could specify \hat{f} using only a portion of the replications, and then assess the generalizability of \hat{f} by seeing how well \hat{f} describes the remaining replications. The more replications we have for testing \hat{f} , the better we can ensure that whatever \hat{f} we use encodes minimal artifacts.

This sets up a tradeoff between descriptive power and robustness. Decomposition (i) lacks the descriptive power to capture across pixel dependencies in pixel color; decomposition (iii) captures features of the observed sequence which may simply be noise. Creating reliable controls hinges on identifying the appropriate decomposition of the data into replicating subunits (see Kruskal, 1988). Statistical inference generalizes reasoning about deterministic patterns to reasoning about *probabilistic* patterns. The fundamental requirement of a pattern—something that repeats—remains unchanged.

No Free Lunch. To summarize, by decomposing (D, θ) into basic building blocks which follow a probabilistic pattern we can create controls (D', θ') from D alone. This strategy turns on an *intestable* assumption: the pattern underlying D persists through θ . That is, the pattern present in the 9 *observed* pixels persists through the 10th and *unobserved* pixel. This assumption is the price we pay for creating *internal* controls (those based on D alone). In contrast, whenever the doctor uses clinical data, she uses *external controls* for Mr. Payne. To see the difference clearly, imagine you want to study yourself. You can experiment on other humans like you (an external source), or you can take some cells from your body (an internal source), experiment on them and extrapolate the conclusions using biological knowledge of how you are (or are not) the sum of your cells—the latter clearly requires more assumptions.

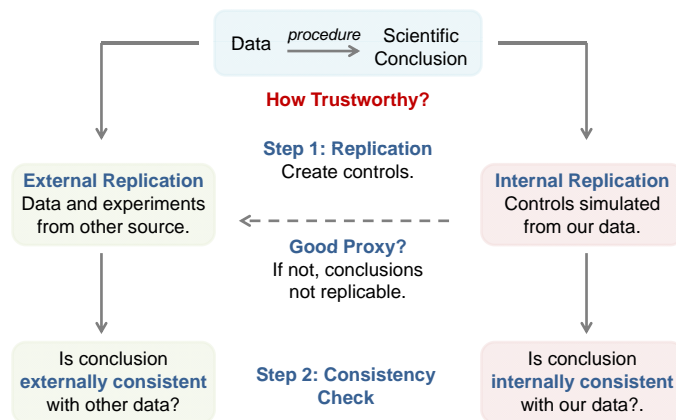


Figure 3: The diagram depicts how one can go about assessing the validity of a scientific conclusion drawn from the data. The gold standard for validation is to replicate the experiment, producing new data, and see whether conclusions drawn from the replicate data are consistent with the current conclusion. Statistical models give us a way to simulate replicate data without having to perform new experiments. If these simulations produce data which closely match those generated by genuine experimental replications, then statistical significance says something meaningful about the real life reproducibility of scientific findings. Otherwise, a statistical uncertainty assessment amounts to only an internal consistency check, i.e., it answers whether the conclusion is reasonable given the current data, not whether the conclusion generalizes.

External replications are a cornerstone of the scientific method. To test the validity of a conclusion, an experiment is repeated over and over to see if consistent conclusions can be

obtained. The *Big Idea* from statistics is the construction of internal controls, giving the data a *built-in capacity* to assess the uncertainty of scientific conclusions drawn from it. That is, having performed this experiment, and before I perform validation experiments, how confident should I be about my conclusions? As discussed above, internal controls are usually of lower quality than external controls because they entail additional assumptions. When uncertainty assessments based on internal replications diverge from those based on external replications, a failure in statistical modelling has occurred, leading to unreplicable research (see Figure 3). Hence the million dollar question is how to create internal replications which closely resemble external replications.

1.3 Frequent Misconceptions: The Meaning of A Statistical Model

This section aims to clarify what we assume *and do not assume* when simulating controls according to a probabilistic pattern/distribution. In the literature, this modeling assumption usually takes the form: “The data *come from* such and such distribution.” This phrasing may give the false impression that...

Misconception 1. A probability model must describe the *generation* of the data.

A more apt description of the model’s job (in inference) is “Such and such probabilistic pattern produces data which *resemble* ours in important ways.” To create replicas (i.e., controls) of the Mona Lisa, one does not need to bring da Vinci back to life—a camera and printer will suffice for most purposes. Of course, knowledge of da Vinci’s painting style will improve the quality of our replicas, just as scientific knowledge of the true data generating process helps us design more meaningful controls. But *for purposes of uncertainty quantification*, our model’s job is to specify a set of controls that *resemble* (D, θ) . Nowhere is this point clearer than in applications involving computer experiments where a probabilistic pattern is used to describe data following a known (but highly complicated) deterministic pattern (Kennedy and O’Hagan, 2001; Conti et al., 2009). We need a *descriptive* model, not necessarily a *generative* model. See Lehmann (1990), Breiman (2001) and Hansen and Yu (2001) for more on this point.

Misconception 2. Because we use a probabilistic model to simulate controls (D', θ') , we must have assumed that D and θ are *random*.

Probability and randomness, so tightly yoked in our minds, are in fact distinct concepts. The language of probability supplies a convenient way to *represent* a set of controls. If we decide that $(D_1, \theta_1), \dots, (D_{10}, \theta_{10})$ act as ten equally compelling controls for (D, θ) , we can mathematically represent this set of controls as a uniform distribution over these ten problems. Nothing has been said about D or θ being random! Randomization inference (see Section 4.2) is a rare instance where the randomness described in the probability model actually corresponds to a physical act of randomization. We can certainly motivate our choice of control problems by conducting thought experiments such as “What if the pixels were generated by a random mechanism?”. But at the end of the day, probability is essentially a tool for bookkeeping, just like the abacus.

1.4 Modeling Complex Datasets by Layering Probabilistic Patterns

In Section 1.2 we decomposed our data into independent replications of a single probabilistic pattern—an i.i.d., *independent and identically distributed*, model. How do we generalize this idea to capture more complicated data structures? An artist often creates a painting through a series of layers—the backmost layer captures broad shapes and general atmosphere (global features) whereas the foremost layers capture fine detail (local features). We can use a similar layering strategy: for us, each “layer” will be a probabilistic pattern, tailored towards capturing data variation at a certain *resolution* level. This is a fundamental idea in many applications, e.g., wavelets (see Daubechies, 2010; Donoho et al., 1995).

To see clearly how the artist analogy applies, suppose we have an image, thought of as an infinite set of pixels, each of which is “area-less”. Our data comprises of the color y_i (a length

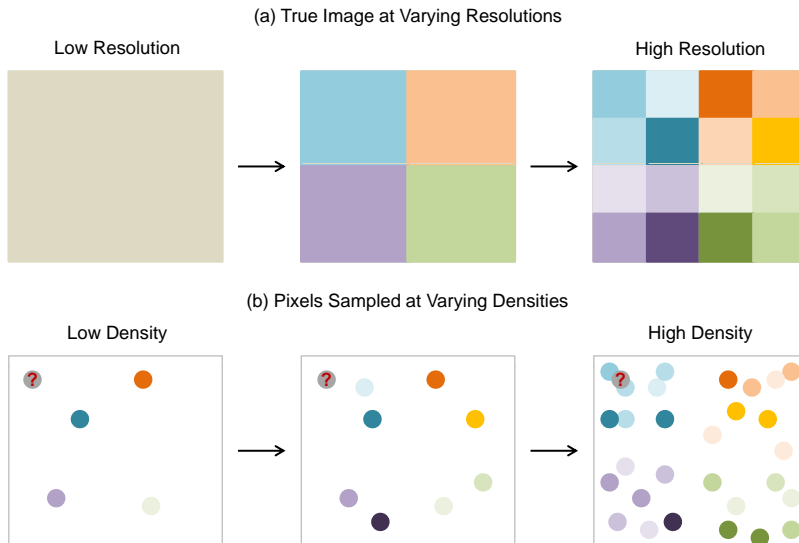


Figure 4: (a) shows the color of the true image averaged over the entire square (left panel), averaged over subsquares (middle panel), or averaged over sub-subsquares (right panel). The true image exhibits smoothness in color variation. (b) shows three datasets (of varying sampling density) comprised of pixels. Each pixel is assumed to be “area-less”.

3 vector specifying the RGB values) and position p_i (the cartesian coordinates) of n randomly sampled pixels; the goal is to impute the color, y , of a pixel at a target location p (see Figure 4b). So $D = \{\{p_i, y_i\}_{i=1}^n, p\}$ and $\theta = y$.

- *A Resolution 1 Model.* Following Section 1.2, we can model color variation at the pixel level using an i.i.d. model, i.e., we assume that the color of each pixel follows some probability distribution *independently* of all other pixels. This model completely ignores the spatial properties of our image: pixels that are closer together are more likely to share a similar color. Control images simulated by this model will exhibit no spatial smoothness. Hence, algorithms which exploit such smoothness properties will perform poorly on our controls even if they perform well on our actual image.

The layering strategy can help us encode smoothness properties into our control images. The model we use below is admittedly naive and has some fatal flaws but its simplicity gets our key points across.

- *A Resolution 2 Model.* Let us divide our image into 4 subsquares as in the mid-panel of Figure 4a. We decompose color variation (across pixels) into two sources: variation in color across the 4 subsquares and variation in color within each subsquare. We assume that the first type of variation follows a probabilistic pattern, $f^{(1)}$, and the second type of variation follows a probabilistic pattern $f^{(2)}$ —this constitutes a resolution 2 model. If in our actual image all pixels in the same subsquare share the same color, then the resolution 2 model would easily pick up on this feature, but the resolution 1 model would not.

Despite being an improvement, the resolution 2 model commits the same sin as the resolution 1 model but just on a finer scale: it will not be able to capture spatial properties of our image *within* each subsquare. Of course, we can further increase our model resolution by subdividing each subsquare into 4 “sub-subsquares” (right panel of Figure 4a) and employ a resolution 3 model. Continuing in this manner, we can build models of arbitrarily high resolution R . Since higher resolution models offer us greater descriptive power, one might be tempted into choosing $R = \infty$. However, we run into the same problem encountered in Section 1.2: descriptive power

comes at the cost of robustness. Highly flexible models tend to “overfit” the data, i.e., our simulated controls will exhibit idiosyncratic features specific to the pixels we happen to observe. Our choice of model resolution, R , should be compatible with the “data resolution.” For example, to meaningfully estimate how smooth our image is, we must sample pixels at sufficiently high density— Figure 4b shows three separate datasets with different sampling densities. If we sampled our image at low density (left panel of Figure 4b), we have no information to assess how smooth the image is within each subsquare.

In summary, “layering” is a widely applicable strategy for modeling complex datasets. A prominent example of its use in statistics is in developing hierarchical models (Gelman and Hill, 2006; Gelman et al., 2013). By combining many layers of probabilistic patterns, we maintain modeling flexibility while retaining the simplest of building blocks—the i.i.d. model—within each layer. As in the basic setting of Section 1.2, a key challenge is balancing a model’s descriptive capacity and its robustness.

Remark 1 *Here, our data are pixels of a fixed image. There is nothing random about the generation of the data. Yet nonetheless, we can describe the data in terms of probabilistic patterns for the purpose of creating control images (D', θ') . This reinforces our point in Section 1.3 that the probabilistic nature of statistical models has little to do with “randomness” in the data generating process.*

1.5 The Weakest Links

This section identifies which components of our model are hardest to specify based on the data alone. In Section 3, we will develop inferential strategies which are robust to these weak links. For concreteness, we first translate the resolution 3 model discussed above into a mathematical form. Let $\mu^{(0)}$ be the “base” color of our actual image (Figure 4a left panel); the superscript indicates that this is layer 0 in our multi-layer model. Let $\left\{ \mu_j^{(1)} \right\}_{j=1}^4$ be the average color over each of the 4 subsquares, $\left\{ \mu_k^{(2)} \right\}_{k=1}^{16}$ be the average color over each of the 16 sub-subsquares, and y_i be the color of a sampled pixel. We decompose pixel color into the following 3 sources of variation:

$$y_i = \mu^{(0)} + \underbrace{\left(\mu_{j(i)}^{(1)} - \mu^{(0)} \right)}_{\text{variation across subsquares}} + \underbrace{\left(\mu_{k(i)}^{(2)} - \mu_{j(i)}^{(1)} \right)}_{\text{variation across sub-subsquares}} + \underbrace{\left(y_i - \mu_{k(i)}^{(2)} \right)}_{\text{variation within sub-subsquares}}.$$

Here $j(i)$ denotes the index of the subsquare and $k(i)$ the index of the sub-subsquare containing pixel i . The resolution 3 model assumes that the base color $\mu^{(0)}$ follows a distribution $f^{(0)}$, the variations $\mu_j^{(1)} - \mu^{(0)}$ (for $j = 1, \dots, 4$) follow a distribution $f^{(1)}$, $\mu_{k(i)}^{(2)} - \mu_{j(i)}^{(1)}$ a distribution $f^{(2)}$, and $y_i - \mu_{k(i)}^{(2)}$ a distribution $f^{(3)}$. How much information is there in the data for specifying the probabilistic patterns $f^{(0)}$, $f^{(1)}$, $f^{(2)}$ and $f^{(3)}$ respectively? Below, we assume that pixels are sampled randomly so that our sample is “representative” of the overall image.

- *The highest resolution pattern (“noise distribution”): $f^{(3)}$.* Suppose we estimate the average color of a sub-subsquare, $\mu_k^{(2)}$, using the average color of the sampled pixels in that sub-subsquare, $\bar{y}_k^{(2)}$. The “residuals” $y_i - \bar{y}_k^{(2)}$ should approximately follow the distribution $f^{(3)}$ (assuming we have sampled enough pixels from each sub-subsquare); if in total we sample n pixels from our image, we essentially observe n repeats of the probabilistic pattern $f^{(3)}$. As a result, it is quite easy to detect misspecification of $f^{(3)}$ using the data. For example, residual plots are often used to diagnose whether a normal distribution is appropriate for describing various aspects of the data.
- *The lowest resolution pattern (“prior distribution”): $f^{(0)}$.* Suppose we estimate $\mu^{(0)}$ using $\bar{y}^{(0)}$, the average color across all sampled pixels, which is an approximate replication of the pattern $f^{(0)}$. Just as we used the residuals $y_i - \bar{y}_k^{(2)}$ to check whether a specification

of $f^{(3)}$ is reasonable, so we can use $\bar{y}^{(0)}$ to assess the fit of $f^{(0)}$ (see Evans and Moshonov, 2006). The key difference of course is that $f^{(0)}$ repeats only once, and hence it cannot be reliably specified based on the data alone.

For medium resolution patterns, there exist no consensus guidelines regarding the number of replications needed for a robust specification. For example, there are 16 across sub-squares differences (corresponding to $\mu_k^{(2)} - \mu_j^{(1)}$) which are approximate replications of $f^{(2)}$. Is this enough to meaningfully specify $f^{(2)}$? See Cox (2006, p.62) and Gelman et al. (2013, p.119) for some discussion. The key takeaway is that in a multi-layer model, the “higher resolution” patterns—which repeat more often throughout the data—are easier to specify than the “lower resolution” patterns.

A natural question then arises: Can we make our inference robust to (ideally, independent of) those patterns which are hard to specify? For example, we might seek a prediction algorithm which performs well on controls (D', θ') regardless of the low resolution pattern $\hat{f}^{(0)}$ used in simulating those controls. This is a primary goal of Frequentist statistics (see Berger, 1985; Efron, 1986). As we shall see, such robustness often comes with a price—loss of *relevance*. The key challenge in statistical inference—as in the doctor’s problem—is then to balance robustness and relevance when creating controls. This topic consumes the remainder of this article; we will study it through an idealized problem detailed below.

Problem Setup: For simplicity, consider data sampled independently from n units, $D = \{y_i\}_{i=1}^n$. We will assume that a simple decomposition of the data suffices in describing its structure: $y_i = g(\theta, \varepsilon_i)$, where g is a known function, θ varies according to a low resolution pattern f_θ and ε_i are independent replications of a high resolution pattern f_ε . For example, $\{y_i\}_{i=1}^n$ may be multiple measurements for a quantity of interest θ with additive error ε_i : $y_i = \theta + \varepsilon_i$. Note that we use θ here as both a parameter in our description of the data and as our quantity of interest because these two often coincide in practice. To perform inference, we need to simulate control problems (D', θ') using specifications \hat{f}_θ and \hat{f}_ε . The wrinkle in our plan is that there is often insufficient information to reliably specify f_θ . Despite the simplicity of this setup, it allows us to access nearly all the important lessons.

Remark 2 *Meaningful specification of f_θ usually requires information external to the data. As a result, f_θ is often called a “prior distribution” (it encodes prior/external knowledge). Despite the special nomenclature, the functional purpose of f_θ and f_ε is the same: to simulate controls. What separates them is the extent to which the data contain information for specifying f_θ versus f_ε . This distinction is a matter of degree. It is counterproductive to treat f_θ and f_ε as categorically different objects since there are an increasing number of datasets which do contain meaningful information for specifying f_θ ; efficient methods for such data require treating f_θ on a more equal footing with f_ε (see Empirical Bayes in Section 4.2)*

2 The Relevance-Robustness Tradeoff

2.1 How Good Are Our Controls?

Let $\hat{\Omega} = \{(D', \theta')\}$ denote our set of simulated control problems, as described in the previous section. The controls in $\hat{\Omega}$ allow us to evaluate the effectiveness of various statistical procedures in the same way that medical procedures are tested through clinical trials. The evaluation process has two basic steps:

1. *Define the evaluation criterion.* We will use Δ' to denote the loss/error of a procedure when applied to a control problem (D', θ') . For example, we may want to test whether a hypothesis H_0 about θ' is true ($T' = 1$) or false ($T' = 0$). A statistical test then takes data as input and outputs an estimate \hat{T}' of the truth value T' . Our test errs ($\Delta' = 1$) whenever $T' \neq \hat{T}'$; otherwise $\Delta' = 0$.

2. Estimate the procedure’s error, Δ , for the target problem, (D, θ) , using a summary of its performance over control problems. For example, we might estimate Δ by the procedure’s average error over $\hat{\Omega}$, denoted $\bar{\Delta}'$. In hypothesis testing, $\bar{\Delta}'$ is the proportion of controls where the test fails, $\hat{T}' \neq T'$.

Figure 5 summarizes the notation and terminology associated with this evaluation process for the three most common inferential tasks: point estimation, set estimation, and hypothesis testing.

Method Type	Error on Actual Problem Δ	Average Error over Relevant Controls $\bar{\Delta}'$	References
Point Estimate <i>Goal:</i> Give our best guess, $\hat{\theta}$, for value of θ .	$L(\theta, \hat{\theta})$ <i>Loss:</i> Specify how “far” θ is from $\hat{\theta}$ via loss function.	<i>Risk:</i> The average loss of an estimator over control problems (D', θ') .	Robinson 1979b Rukhin 1988, Lu and Berger 1989, Fourdrinier and Wells 2012
Set Estimate <i>Goal:</i> Identify set, $C(D)$, of likely values for θ .	$I(\theta \notin C(D))$ * <i>Coverage:</i> Does our set contain the true value of θ ?	<i>Non-Coverage Probability:</i> Proportion of times a set estimate, e.g. interval estimate, fails to contain the true value of θ' .	Casella 1992, Goutis and Casella 1995, Robinson 1979a, Berger 1988
Hypothesis Test <i>Goal:</i> Should we reject a null hypothesis, H_0 , based on evidence from data?	$I(\hat{T} \neq T)$ <i>Type I or II Error:</i> Do we falsely reject or falsely accept H_0 ?	<i>Error Probability:</i> The test’s rates of false rejection and false acceptance when applied to control problems.	Hwang et al. 1992, Berger et al. 1994, Berger 2003

* $I(\text{statement})$ denotes the indicator function: it equals 1 if the statement in parentheses is true and 0 otherwise.

Figure 5: Three primary ways to summarize information about θ are point estimates, set estimates, and hypothesis tests. For a summary to be meaningful, we need to assess its quality when applied to the target problem (D, θ) , i.e., we need to estimate Δ . We do so by evaluating the performance of each procedure over control problems. The references listed are a small sample of articles discussing how the relevance of our control problems affects the accuracy of the resulting estimates $\bar{\Delta}'$ of Δ .

A natural question arises at this juncture: How well does a procedure’s average performance over controls predict its actual performance Δ (see Kiefer, 1976, 1977; Brown, 1978; Sundberg, 2003)? For example, suppose a researcher finds a formula for a “95% confidence interval” and uses it to produce an interval estimate for a quantity θ . Is 95% a good estimate for whether his interval actually contains θ ? The answer obviously depends on the control problems used to compute the “95%” figure.

As foreshadowed by our medical analogy, the initial set of controls, $\hat{\Omega}$, usually contains problems that are a poor imitation of (D, θ) . We can usually select a subset $\hat{\Omega}_{\text{rele}} \subset \hat{\Omega}$ of controls that more closely *match* the observed features of (D, θ) , hence is more *relevant*. This idea has a long history: Fisher (1959) described it as finding a “reference set” for our problem and proposed the criterion of “recognizable subsets” (see Zabell, 1992, Section 7.2) to judge whether more matching is needed. While conceptually simple, we quickly run into trouble when hammering out the details of the matching process. In particular, the relevance of a control problem turns out to be highly sensitive to the specification of \hat{f}_θ in our model (Section 2.3)—which we know from Section 1.5 can be quite unreliable. Since gains in relevance may be washed out by loss of robustness, we need to strike the right balance. For example, the reader may be aware that there are 95% “Frequentist” confidence intervals and 95% “Bayesian” confidence intervals—both have a success rate of 95% but with respect to *different* choices of control problems. In particular, the Bayesian favors relevance and matches (much) more than the Frequentist. Which of the two made the better tradeoff? To give a mathematically satisfactory answer to this question is quite difficult; the interested reader should see the references in Figure 5. It turns out, however, that by returning to our medical analogy and studying how a doctor

selects controls for Mr. Payne, we obtain a much more tractable problem which holds the same lessons as the original. Readers satisfied with the heuristics given above may skip the following section, which is inevitably more technical.

2.2 A Cost Benefit Analysis of Matching

The doctor wants to know y , Mr. Payne’s response if he were assigned the newly approved treatment B . Just as we estimate Δ by $\bar{\Delta}'$, the average error over control problems, so the doctor estimates y using \bar{y}' , the average response to treatment B over some set of control patients. We can hopefully improve upon this initial estimate by exploiting the fact that some patients are more relevant for Mr. Payne than others. But achieving higher relevance comes with the obvious cost of a reduction in the number of available controls. To what degree then should controls match the individual of interest?

The ideal case of infinitely many matches. We begin our analysis by removing one side of the tradeoff: the constraint that our clinical trial enrolls only a finite number of patients. Let Ω denote an infinite set of control patients and $\Omega_{\bar{a}}$ the subset with background variables $\bar{x}' = \bar{a}$. In particular, $\Omega_{\bar{x}}$ consists of those patients who match Mr. Payne’s observed characteristics completely. We show that the average treatment response for patients in $\Omega_{\bar{x}}$, denoted $\mu_{\bar{x}}$, is a better estimate of Mr. Payne’s response y than the average response for patients in Ω , denoted μ . Our criterion will be the *mean-squared error* (MSE). Ideally we want to compare $(y - \mu_{\bar{x}})^2$ and $(y - \mu)^2$, but because they are unknown, we estimate them by their corresponding averages over Ω (or $\Omega_{\bar{x}}$). That is, for each (\bar{x}', y') in Ω , we compute the accuracy of μ in predicting y' ; we then repeat this assessment for $\mu_{\bar{x}'}$. Because $(y' - \mu)^2 = (y' - \mu_{\bar{x}'})^2 + (\mu_{\bar{x}'} - \mu)^2 + 2(y' - \mu_{\bar{x}'})(\mu_{\bar{x}'} - \mu)$, and the last term is zero when averaged over Ω , we have

$$\text{Ave}_{\Omega} [(y - \mu)^2] - \text{Ave}_{\Omega} [(y - \mu_{\bar{x}'})^2] = \text{Ave}_{\Omega} [(\mu_{\bar{x}'} - \mu)^2]. \quad (1)$$

Here, $\text{Ave}_{\Omega} [\cdot]$ is the average of the bracketed expression across controls in Ω . We see that the error from using $\mu_{\bar{x}}$ to predict y is (on average) no worse than that from using μ . And the expected improvement— $\text{Ave}_{\Omega} [(\mu - \mu_{\bar{x}'})^2]$ —is precisely the average (squared) difference between the more relevant and less relevant estimates of y' .

Finite-sample complications. Real clinical trials have only a finite number of patients. Let $\hat{\Omega}$ denote this finite set and $\hat{\Omega}_{\bar{a}}$ the subset of patients with $\bar{x}' = \bar{a}$; analogously, $\hat{\mu}$ and $\hat{\mu}_{\bar{a}}$ denote the average patient response in $\hat{\Omega}$ and $\hat{\Omega}_{\bar{a}}$ respectively. In the case of unlimited controls, on average $\mu_{\bar{x}}$ is a better predictor of y than μ ; with only a finite number of controls, it is now possible for $\hat{\mu}$ to outperform $\hat{\mu}_{\bar{x}}$. This is because we are implicitly engaging in a two-stage estimation process—(a) estimate y by μ (or $\mu_{\bar{x}}$) and (b) estimate μ by $\hat{\mu}$ (or $\mu_{\bar{x}}$ by $\hat{\mu}_{\bar{x}}$). Since few individuals match Mr. Payne’s background completely, it is harder to estimate $\mu_{\bar{x}}$ than μ . This increase in estimation error can wipe out benefits from matching.

Finding the sweet spot. When the cost of estimating $\mu_{\bar{x}}$ by $\hat{\mu}_{\bar{x}}$ is too high, we can compromise by using controls who match Mr. Payne with respect to some, but not all, variables. Let $\bar{x}'_{(0)}, \bar{x}'_{(1)}, \dots, \bar{x}'_{(R)}$ denote sequentially expanding collections of background variables: for example, $\bar{x}'_{(1)}$ could be age, $\bar{x}'_{(2)}$ age *and* blood pressure, etc. We use $\bar{x}'_{(0)}$ to denote no matching and $\bar{x}' = \bar{x}'_{(R)}$ to denote matching on all the measured variables. A natural question then is: *given such a sequence*, for what value of r between 0 and R , does the incremental increase in estimation error exceed the incremental gain from matching? (An even harder question is how to order the sequence so that the “most important” predictors appear first; see Meng 2014 for a discussion.)

To answer this question, we want to compare $(y - \hat{\mu}_{\bar{x}_{(r)}})^2$ versus $(y - \hat{\mu}_{\bar{x}_{(r+1)}})^2$, i.e., matching on r versus $r + 1$ background variables. Again, these two prediction errors are unknown, so we will estimate them by their corresponding averages over control patients, say in Ω . The average prediction error for $\hat{\mu}_{\bar{x}'_{(r)}}$ minus the average prediction error for $\hat{\mu}_{\bar{x}'_{(r+1)}}$ decomposes as follows:

$$\underbrace{\text{Ave}_{\Omega} [(\mu_{\bar{x}'_{(r)}} - \mu_{\bar{x}'_{(r+1)}})^2]}_{\text{gain in relevance}} - \underbrace{\text{Ave}_{\Omega} [(\hat{\mu}_{\bar{x}'_{(r+1)}} - \mu_{\bar{x}'_{(r+1)}})^2] - \text{Ave}_{\Omega} [(\hat{\mu}_{\bar{x}'_{(r)}} - \mu_{\bar{x}'_{(r)}})^2]}_{\text{loss in robustness}}. \quad (2)$$

The second term compares how difficult it is to estimate $\mu_{\bar{x}'_{(r+1)}}$ versus $\mu_{\bar{x}'_{(r)}}$ —since it is harder to find controls which match on $r + 1$ variables as compared to r , $\mu_{\bar{x}'_{(r+1)}}$ is usually harder to estimate. So, the second term is usually positive; it represents the loss in robustness from an additional step of matching. The first term is analogous to the RHS of equation (1), i.e., it represents the benefit of matching. Therefore, expression (2) lays out mathematically the relevance robustness tradeoff.

How this all relates to statistical inference. It may be difficult at first to see why the lessons above transfer over to statistical inference. After all, what limits the doctor’s ability to match is the *finite* nature of her clinical trial, whereas statisticians can simulate as many control problems (D', θ') as they want, at least in principle. The key is this: having a finite clinical trial means that if we rerun the trial by sampling n new patients, our inferences would change. That is, $\hat{\Omega}$ is unstable. The ideal trial Ω —by virtue of enrolling “everyone”—is free of the idiosyncrasies of any sample; it is stable. In the statistical context, the set of control problems, $\hat{\Omega}$, is also “unstable” because we cannot specify the patterns \hat{f}_θ and \hat{f}_ε with perfect certainty. In both contexts, the instability of $\hat{\Omega}$ is magnified when considering subsets of controls—in this way, matching erodes robustness.

2.3 Matching: The Doctor Does It, So Should the Statistician

The more closely a clinical trial patient matches Mr. Payne’s observed characteristics, \bar{x} , the more relevant the patient is. Similarly, relevant control problems, (D', θ') , will match our target problem, (D, θ) , with respect to its “observed characteristics”: D . Complete matches are controls with $D' = D$. This correspondence between \bar{x} and D makes intuitive sense: if the data in a control problem fail to mimic important features of our actual data, evaluating a procedure on the former tells us little, or may even mislead us, about its performance on the latter. In the words of Kiefer (1976, 1977), we need to distinguish those problems with “lucky” data from those with “unlucky” data, as illustrated below.

Example 2 (Data Precision; Cox 1958). Suppose we have n independent measurements on the weight of an object. The sample size, n , is a feature of the data (an often forgotten fact); since problems with very different sample sizes are not comparable, we want our controls to have a matching sample size of n . Suppose we also know that the task of taking the i th measurement was randomly assigned to one of two labs; the first lab produces very precise estimates, the second lab noisier ones. By chance, all but two measurements in *our* sample came from Lab 2. The lab assignments are also an important feature of the observed data. If a control problem has data where the majority of measurements originate from the more precise Lab 1, an estimator’s effectiveness on the control may be an overly optimistic assessment of its error when applied to our data. Hence, we should also match on lab assignment.

Sample size and lab assignment help determine the *precision* of our data; in other contexts, other features of the data may help determine precision. For example, suppose we are analyzing data from a randomized clinical trial comparing treatments A and B . By chance, 80% of females received treatment A ; this imbalance makes it harder to disentangle the treatment effect from any gender effect. If control datasets do not preserve this gender imbalance, the estimate of the treatment effect would appear to be much more precise than it actually is, causing us to be overconfident. Hence Cox (1982) and Rosenbaum (1984) suggest keeping only those control problems where the covariate balance is sufficiently similar to the observed balance. Both examples above teach us that the precision of control data should match that of our actual data. This seemingly obvious principle is violated surprisingly often in practice; see Fraser (2004) for discussion.

In problems involving *independent* measurements, matching on the sample size is standard practice. The qualitative effect of matching—larger sample sizes imply greater precision, smaller ones less precision—is robust to our initial selection of control problems, $\hat{\Omega}$, i.e., robust to model

misspecification. As Example 3 below shows, this robustness property is not shared by all features of the data. The decision whether to match on these other features is much more controversial. In some situations, even matching on the sample size can lead to non-robust inferences (see Rosenbaum and Rubin, 1984).

Example 3. Let y be a single measurement for the weight, θ , of an object; the natural estimator of θ is y itself. To assess its accuracy, we simulate controls as follows: $y' = \theta' + \varepsilon'$, where the measurement errors ε' vary according to a distribution f_ε (supplied by the scale manufacturer) and we set $\theta' = c$. For each choice of the simulation parameter c , we obtain a different set of control problems. However, our estimator’s average error $|y' - \theta'|$ over $\hat{\Omega}$ turns out to be independent of c . This robustness vanishes when we match on the measurement itself; the average error over controls with $y' = y$ is $|y - c|$, which can range anywhere from 0 to ∞ depending on c . Intuitively, the value of the measurement y tells us nothing about its precision *unless* we have some external information about θ . If we do have this information, say $\theta < 10$, then the value of y can be very informative about its precision, e.g., $y > 10$ would be an imprecise measurement.

In order to justify matching on y in this problem (as is done in Bayesian inference), we must know something (reliable) about θ . When we do not, can we somehow specify an “objective/non-informative” choice of \hat{f}_θ and obtain a sensible inference from matching on y ? This quest underlies objective Bayesian analyses (see Kass and Wasserman, 1996; Berger, 2006). Such analyses usually choose \hat{f}_θ with the goal of harmonizing the inference from matching on y with the inference from not matching on y . Thus, while we nominally obtain an “individualized” inference, the matching has no value added. In this example, there is no way to obtain *meaningful* “individualization” without genuine information about θ .

The features of the data described in Example 2 are examples of *ancillary* statistics (see Fraser, 2004; Ghosh et al., 2010); the idea of matching on ancillaries dates back to Fisher (1925, 1934). The name, ancillary, comes from the fact that such features of the data can be simulated *without* specifying \hat{f}_θ (this gives an intuitive way to see why matching on such features does not erode the robustness of our inference). In contrast, the measurement y' in Example 3 is not ancillary—to simulate y' , we first simulate θ' . Ancillary statistics tell us something about the precision of our data even when we know nothing about θ . Non-ancillary statistics also carry information on the precision of our data, but the effective use of this information often requires prior knowledge about θ . Many features of the data lie in between the extremes of Example 2 and Example 3, i.e., the information they contain about precision is partially sensitive to \hat{f}_θ (see Examples 4 and 5). The relevance robustness tradeoff becomes much harder to navigate in such cases.

Example 4 (Accounting for Selection Bias). Using data from n patients, we wish to identify genetic markers which influence an individual’s cholesterol level. We estimated each marker’s effect by looking at the difference in average cholesterol between those with and without the marker; denote our estimates by $\hat{\theta}_m$ (for $m = 1, \dots, M$). We also computed a p-value for each marker assessing whether $\hat{\theta}_m$ is significantly different from 0. Only markers with p-value below some threshold are selected for further study. But how accurate are the estimates $\hat{\theta}_m$ for the *selected* markers?

Let us focus on marker 1. The key is to realize that there are two types of datasets—those in which marker 1 is selected and those in which it is not—and that the accuracy of $\hat{\theta}_m$ differs (often dramatically) between the two types. To see this, note that two factors help determine whether a marker is selected: (i) the magnitude of the marker’s impact on cholesterol and (ii) luck/chance. Luck in this case means that the magnitude of our estimate $|\hat{\theta}_m|$ is upward biased compared to the true magnitude (i.e., luck helps marker m get selected). When we perform selection, we are targeting not only markers with large effects, but also markers for which the current dataset is “lucky”. Our original estimates for the selected markers likely overstate their impact. Bias induced by selection is commonly referred to as “winner’s curse” (Ioannidis et al., 2001). We want to remove this curse.

Suppose we have some model for simulating control problems (D', θ') . The reasoning above tells us that we should restrict attention to controls where the same markers are selected as in our actual study; without this matching step, our inference ignores selection bias. After matching, we can compute the bias of $\hat{\theta}'_m$ over the relevant controls; this acts as an estimate of the actual bias of $\hat{\theta}_m$. However, there is a problem: the bias of $\hat{\theta}'_m$ over the relevant controls is sensitive to how we simulate θ' (i.e., sensitive to \hat{f}_θ). If a marker's true effect size, $|\theta'_m|$, is large, it needs very little luck to be selected, so the bias of $\hat{\theta}'_m$ will be small; the converse holds when $|\theta'_m|$ is small. To effectively estimate the winner's curse by matching, we need a reliable specification of \hat{f}_θ ; fortunately this is often possible when the number of markers is large (see Efron, 2010).

Remark 3 *The above example is a variation on the statistical practice of model selection. When fitting a model of outcome on predictors, a selection step may be introduced to screen out unimportant predictors, e.g., using criterion such as AIC and BIC (see Burnham and Anderson, 2002). This practice is becoming increasingly common because modern applications often involve so-called “large-p small-n” data where the number of predictors exceeds the number of observations. One may then wish to estimate and build confidence intervals for the parameters in the selected model. Ignoring the selection bias will lead to seriously flawed inferences (see Leeb and Pötscher, 2005).*

Example 5 (The Strength of Evidence). This example comes from Dempster (1997) and is extended using ideas in Berger (2003). We want to test a null hypothesis, $H_0 : \theta = 0$, against an alternative, $H_1 : \theta = 1$. Suppose we observe a p-value $p = 0.049$, and reject the null hypothesis according to a pre-determined testing procedure: reject H_0 if $p \leq 0.05$. We want to assess the probability that we have made a false rejection. So we evaluate our testing protocol over some relevant subset of control problems (p', θ') in $\hat{\Omega}$ (for simplicity, we assume here that p' preserves all the information in the data). The question is: which subset?

Case 1. *Use All Control Problems in $\hat{\Omega}$:* p-values are constructed to follow a uniform distribution when $H_0 : \theta' = 0$ is true. Hence, our error rate (proportion of times $p' \leq 0.05$) over control problems with $\theta' = 0$ is 5%; this is the Type I error of our test. Suppose the p-value over problems with $\theta' = 1$ follows a Beta(0.02, 1.35) distribution (see Figure 6a)—that is, we are more likely to see very small p-values when H_1 is true. Under this Beta distribution, 5% of datasets with $\theta' = 1$ have $p' > 0.05$ so our error rate under H_1 (Type II error) is also 5%. Thus, *regardless of \hat{f}_θ* , the error rate of our test over all problems in $\hat{\Omega}$ is 5%.

Despite its appealing robustness, we notice some unsettling facts about the above inference.

- If we had observed $p = 10^{-8}$, our error assessment (using $\hat{\Omega}$) would still be 5%. Yet the chance of false rejection is intuitively much smaller when $p = 10^{-8}$ than when $p = 0.049$.
- In Figure 6a, the density of the p-value distribution at $p = 0.049$ is higher when H_0 is true than when H_1 is true. The ratio of the density under H_1 to that under H_0 is $LR = 0.38$. We call this the *likelihood ratio*. A $LR < 1$ says that the data supports H_0 more than H_1 ! Given this, a 5% error rate seems an inaccurate assessment of our uncertainty.

These incongruities have a simple explanation. The error rate over $\hat{\Omega}$, measures the *average accuracy* of a control problem. But a rejection based on $p' = 10^{-8}$ is much more convincing than one based on $p' = 0.049$. So the average accuracy will underestimate the actual accuracy when $p = 10^{-8}$ and overestimate it when $p = 0.049$.

To fix this, we need to match on some measure of precision. The likelihood ratio seems to capture precision reasonably well, but we will make one small modification to it. Note that a testing problem with $LR' = \frac{1}{2}$ (evidence favors H_0 by factor of 2) seems as difficult as one with $LR' = 2$ (evidence favors H_1 by factor of 2). This symmetry motivates us to define as *relevant* those controls for which either $LR' \approx LR$ or $LR' \approx LR^{-1}$; more compactly, controls such that $|\log LR'| - |\log LR| \approx 0$ (see Berger et al., 1994, for an argument which leads to matching on ancillary features).

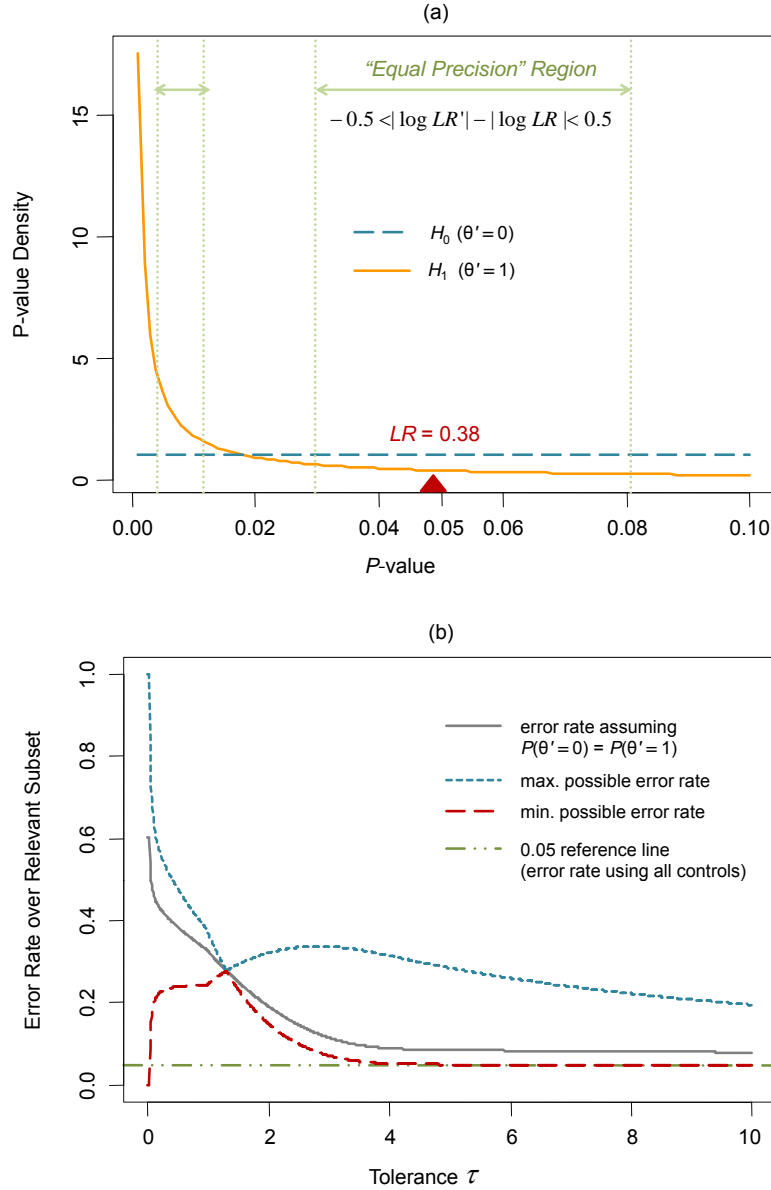


Figure 6: Plot (a) depicts the densities of the p-value distribution when H_0 or H_1 are true. The density under H_1 is higher for small p , i.e., small p-values favor H_1 . The red triangle marks the position of the observed p-value, whose density under H_0 exceeds its density under H_1 (the evidence favors H_0). The dashed vertical lines demarcate an “equal precision region” (this is a disjoint region): only controls with an observed p-value, p' , lying in this region are considered sufficiently relevant for the actual problem. Here, “sufficiently relevant” is specified through the tolerance $\tau = 0.5$. The solid gray line in (b) plots the error rate of our test over $\hat{\Omega}_{\text{rele}}$ for various values of τ (when \hat{f}_θ assigns equal probability to H_0 and H_1). The dashed lines show the maximal and minimal possible error rates over $\hat{\Omega}_{\text{rele}}$ when \hat{f}_θ is allowed to be arbitrary. Sensitivity to \hat{f}_θ increases as $\tau \rightarrow 0$, that is, as we increase the relevance.

Case 2. $\hat{\Omega}_{\text{rele}}$ contains those problems (p', θ') in $\hat{\Omega}$ such that $-\tau \leq |\log LR'| - |\log LR| \leq \tau$. The tolerance τ measures the degree to which we require our controls to match our actual problem; for $\tau = 0.5$, $\hat{\Omega}_{\text{rele}}$ includes only controls with observed p-value, p' , in the “equal precision region” (see Figure 6a). When $\tau = 0$, $LR' = 0.38$, or $1/0.38$ —our control has the same precision as the actual problem. When $\tau = \infty$, we use all controls in $\hat{\Omega}$; this is equivalent to Case 1.

Unlike Case 1, the error rate of our test over $\hat{\Omega}_{\text{rele}}$ depends on our specification of \hat{f}_θ ; Figure 6b plots the maximal and minimal error rates (dashed lines) that can be achieved by varying \hat{f}_θ (this helps us assess how sensitive our inference is). Figure 6b also plots the error rate over $\hat{\Omega}_{\text{rele}}$ when \hat{f}_θ assigns equal probability to H_0 and H_1 (a popular choice).

As we increase the relevance ($\tau \rightarrow 0$), inferences drawn from $\hat{\Omega}_{\text{rele}}$ becomes less and less robust: the error rate ranges anywhere from 0 to 1 depending on how we specify \hat{f}_θ . However, for almost all reasonable choices of τ , the range of plausible error rates lies above 5%—the error rate when we do not match. This suggests that, at the least, the qualitative impact of matching—to correct the downward bias of our initial error estimate—is quite robust even if the degree of adjustment is not. Intuitively, the more *heterogeneous* the problems in $\hat{\Omega}$, the more important it is to match. In this example, the performance of our testing procedure is sufficiently different when $p' = 10^{-8}$ than when $p' = 0.049$ to justify some degree of matching—of course, deciding the exact degree is a much harder problem! Interestingly, Berger and Wolpert (1988, Example 4b) give an even more extreme form of this example where $\hat{\Omega}$ is so heterogeneous that complete matching is warranted.

In summary, Example 2 is the poster child for matching (relevance), Example 3 the burning beacon of its potential dangers; Examples 4 and 5 fall somewhere in between. In choosing the right relevance robustness tradeoff for a problem, it is helpful to think about how it maps in terms of these canonical examples.

3 Robustness Through Symmetry and Invariance

To increase the relevance of controls in our approximating population, we match on features of the data. Is there any way to increase robustness? More precisely, can we design procedures which perform well regardless of \hat{f}_θ ? Consider Example 5: the average error of our test over $\hat{\Omega}$ was 5% for any choice of \hat{f}_θ . This robustness derives from the *symmetry* or *invariance* of the procedure. The accuracy of our hypothesis test (over $\hat{\Omega}$) is invariant to whether the null or alternative hypothesis is true. As the following example demonstrates, *conditioning* (the technical term for matching) on non-ancillary features of the data breaks the symmetry of a procedure. The implication is that it is hard to construct procedures which are symmetric not only over $\hat{\Omega}$ but also over *subsets* of $\hat{\Omega}$. The intuition for this phenomenon is simple: it’s hard designing medicine that works equally well on all individuals regardless of background.

Example 6 (Too High a Price for Relevance). A doctor wants to predict your disease status: sick ($\theta = 1$) or healthy ($\theta = 0$)? A test for the disease yields a positive result, $y = 1$. The company producing the test tells the doctor that, in a clinical trial $\hat{\Omega}$, sick individuals tested positive with 90% probability while healthy ones tested positive with 10% probability. To decide whether to initiate treatment, the doctor asks for the probability that *your test result* is wrong.

Observe that the test is wrong 10% of the time for sick individuals (false negative) and 10% of the time for healthy individuals (false positive), i.e., the test’s accuracy is invariant to the patient’s underlying disease status. Hence, the test’s error rate with respect to the clinical trial population, $\hat{\Omega}$, is 10%. This assessment does not depend on knowing the proportion of sick individuals ($\theta' = 1$) in the clinical trial. This proportion, which is potentially unknown to the doctor, acts as \hat{f}_θ in our example.

Not all clinical trial patients are relevant *for you*: your test came back positive while the clinical trial contains individuals with negative results. What the doctor actually seeks is the accuracy of the test among the subset, $\hat{\Omega}_{\text{rele}}$, of individuals *who tested positive*. Among this

subset of patients, the accuracy of the test is no longer invariant to the patient’s underlying disease status: the test is 100% accurate for the sick and 0% accurate for the healthy. The test is symmetric over $\hat{\Omega}$ but not over $\hat{\Omega}_{\text{rele}}$. Due to this loss of invariance, the test’s accuracy among patients in $\hat{\Omega}_{\text{rele}}$ becomes sensitive to the proportion of sick individuals. In fact, the error rate can range anywhere from 0% to 100%!

Example 6 shows how “symmetry” allows us to make inferences that are robust to misspecification of \hat{f}_θ (see Rubin, 1978, for an application of this idea in the case of randomized experiments). The price, of course, is that we need to limit the degree to which we “individualize” our inference—procedures which act symmetrically over a set of controls often behave asymmetrically over its subsets. The pressing question now is how to construct symmetric procedures. It turns out that the study of symmetry is best initiated by studying *asymmetry*.

3.1 Asymmetric Procedures and the Plug-in Principle

A procedure may behave symmetrically in one statistical problem and asymmetrically in another. To illustrate this point, we consider a confidence interval which behaves asymmetrically because it “violates” the problem’s natural symmetries.

Example 7 (An Asymmetric Procedure). For quality control purposes, a company records failure times y_i for n randomly selected batteries. The goal is to find a lower bound on the average time to battery failure, θ . Clearly, we cannot bound θ deterministically, except for the useless bound zero. To create a probabilistic lower bound we begin with an estimate of θ , e.g., the sample mean \bar{y}_n . A seemingly innocuous construction is to then choose a buffer b so that the lower bound $\bar{y}'_n - b$ lies below θ' for a desired fraction of controls; here $D' = \{y'_i\}_{i=1}^n$.

To simulate controls, the company uses a multiplicative model for the failure times: $y'_i = \theta' \varepsilon'_i$, where the mean failure time θ' is simulated from some \hat{f}_θ , and the i th battery’s deviation from the mean, $\varepsilon'_i > 0$, is independently drawn from a distribution \hat{f}_ε with mean 1. The difficulty lies in specifying a reasonable \hat{f}_θ , which requires prior knowledge of the average failure time. For our inference to be robust, we want to choose b such that $\bar{y}'_n - b$ is a good lower bound regardless of \hat{f}_θ . Does a non-trivial choice exist?

The answer is no. We have purposely cooked up an “incompatibility” in our example: our bound $\bar{y}_n - b$ implicitly assumes an additive structure for battery heterogeneity, but our controls are generated from a multiplicative model. This incompatibility makes the performance of our bound ultra-sensitive to any choice of \hat{f}_θ . Since $\bar{y}'_n = \theta' \bar{\varepsilon}'_n$, the lower bound succeeds ($\theta' > \bar{y}'_n - b$) for a control problem if and only if $\bar{\varepsilon}'_n < 1 + b/\theta'$. We can think of $1 + b/\theta'$ as an upper bound for $\bar{\varepsilon}'_n$ with buffer b/θ' , which gets arbitrarily close to zero as we increase θ' . The lower bound’s effectiveness therefore is asymmetric, i.e., it depends on the particular value of θ' , yielding no useful lower bound that will work well for all \hat{f}_θ .

Example 7 shows that we cannot assess the effectiveness of an asymmetric procedure in a manner robust to \hat{f}_θ . Hence, there is little reliable information for deciding whether to use that procedure to analyze our dataset—if we use it, we are at the mercy of luck or unluck. Yet asymmetric statistical procedures are often unavoidable simply because non-trivial symmetric procedures do not exist. A crude but convenient way (see Efron, 1998) to proceed is to “plug-in” an estimate of θ when generating controls, i.e., set $\theta' = \hat{\theta}$ where $\hat{\theta}$ is our estimate. This is equivalent to specifying \hat{f}_θ as a point mass at $\hat{\theta}$; the *parametric bootstrap* (Efron and Tibshirani, 1994; Davison and Hinkley, 1997) and *Empirical Bayes* (see Section 4.2) are common applications of this strategy. The plug-in approach yokes our uncertainty assessment to an initial estimate of θ . If a procedure’s effectiveness is highly sensitive to the value of θ (i.e., highly asymmetric), the resulting inference may be worthless for even small errors in our initial estimate. That is, our error assessment will be least reliable precisely when our estimation error is greatest—a most unattractive quality! The plug-in approach therefore works best when our procedure is nearly symmetric (this point adds additional motivation for understanding how to create symmetric or nearly symmetric procedures). Despite its dangers, plugging in is popular in practice because of its ease of implementation, but it should really be treated as a weapon of desperation.

3.2 Symmetry/Invariance Creation Via Pivots and Minimaxy

We now discuss two strategies for constructing procedures that are symmetric in θ or approximately so. Such procedures have performance guarantees that are robust to the user’s specification of \hat{f}_θ . We may desire robustness to \hat{f}_ε as well as \hat{f}_θ . Extensions of the strategies below can often lead to procedures that are symmetric in both ε and θ ; the interested reader can follow the path to randomization inference (see Figure 8). Of course, the cost of such robustness is a further drop in relevance (see Basu, 1980).

Creating Symmetry/Invariance via Pivots. As motivation, consider this observation from Example 7: the representation $\bar{y}'_n = \theta' \cdot \bar{\varepsilon}'_n$ induces a duality between θ' and $\bar{\varepsilon}'_n$ —a lower bound ($\bar{y}'_n - b$) for θ' induces an upper bound ($1 + b/\theta'$) for $\bar{\varepsilon}'_n$ and vice versa. The fact that the effectiveness of this procedure depends on θ' is self-evident by considering the upper bound for $\bar{\varepsilon}'_n$ (which is small if θ' is large). Intuitively, we can “symmetrize” our procedure by using an upper bound for $\bar{\varepsilon}'_n$ that is *independent* of θ' . The performance of our procedure will then be *invariant* to the choice of \hat{f}_θ . For example, we can choose a constant c such that $\bar{\varepsilon}'_n < c$ for 95% of the controls in $\hat{\Omega}$. By duality, we automatically obtain $\theta' > \bar{y}'_n/c$ for 95% of the controls in $\hat{\Omega}$. Note that:

- The distribution of $\bar{\varepsilon}'_n$ depends only on \hat{f}_ε so our choice of c does not depend on \hat{f}_θ . As a result, the guarantee— $\theta' > \bar{y}'_n/c$ for 95% of controls in $\hat{\Omega}$ —holds for any choice of \hat{f}_θ .
- A lower bound of the form \bar{y}'_n/c reflects the multiplicative model used to simulate controls; in contrast a lower bound of the form $\bar{y}'_n - b$ assumes an incompatible additive model.

By switching our attention from θ' to $\bar{\varepsilon}'_n$, we decouple the performance of our procedure from the choice \hat{f}_θ . That is, we achieve symmetry/invariance. The duality between θ and $\bar{\varepsilon}_n$ —and the switching strategy it allows—lies at the heart of *Fiducial inference* (Fisher, 1930, 1935), which led to the development of confidence intervals in Neyman (1934). The general recipe may be summarized as follows. Identify a quantity $t(\varepsilon')$ with the following properties (in Example 7, $t(\varepsilon') = \bar{\varepsilon}'_n$):

- $t(\varepsilon')$ is dual to θ' , i.e., inference for $t(\varepsilon')$ is equivalent to inference for θ' . Rather than design procedures to estimate θ' , we can design procedures to estimate $t(\varepsilon')$.
- The distribution of $t(\varepsilon')$ over $\hat{\Omega}$ depends only on \hat{f}_ε and *not* on \hat{f}_θ ; such a quantity is said to be *pivotal*. It is usually simple to find “estimators” of pivotal quantities whose effectiveness is decoupled from our choice of \hat{f}_θ . In our example, the value of c which makes $\bar{\varepsilon}'_n < c$ for 95% of controls depends only on \hat{f}_ε .

This approach gives us a systematic way of constructing symmetric/invariant procedures (see Fraser, 1968; Dawid and Stone, 1982; Hannig, 2009).

Pivots and the Relevance Robustness Tradeoff. In principle, we can replace $\hat{\Omega}$ in Criterion II with $\hat{\Omega}_{\text{rele}}$; this lets us create procedures that are symmetric over *subsets* of controls. However, we encountered a key phenomenon in Example 6: matching breaks symmetry/invariance. This implies that the more matching we do, the harder it will be to (non-trivially) satisfy Criteria I and II. Again, consider Example 7.

Case 1: We perform no matching: $\hat{\Omega}_{\text{rele}} = \hat{\Omega}$. The distribution of $\bar{\varepsilon}'_n$ over $\hat{\Omega}$ depends only on \hat{f}_ε and not \hat{f}_θ ; hence $\bar{\varepsilon}'_n$ is pivotal.

Case 2: We use only completely matched controls, i.e., $y'_i = y_i$ for all i . The constraint $\bar{y}_n = \bar{y}'_n = \theta' \bar{\varepsilon}'_n$ induces dependence between θ' and $\bar{\varepsilon}'_n$ over the set of matched controls. The distribution of $\bar{\varepsilon}'_n$ among controls with $\bar{y}'_n = \bar{y}_n$ now depends on \hat{f}_θ , i.e., $\bar{\varepsilon}'_n$ is no longer pivotal *with respect to* $\hat{\Omega}_{\text{rele}}$. Our usage of the term “pivotal” here is an extension of the standard definition given above.

Whereas it is straightforward to find procedures with robust performance over $\hat{\Omega}$, it is difficult to achieve similar robustness over subsets of $\hat{\Omega}$. Most often, only “absurd” estimators can guarantee such robustness. For example, if $\hat{\Omega}_{\text{rele}}$ is the set of complete matches in our current example, only two (non-randomized) interval estimates exist whose performance over $\hat{\Omega}_{\text{rele}}$ is

insensitive to \hat{f}_θ : the interval, $(0, \infty)$, which is 100% correct and 100% unhelpful, or the empty set, which is clearly useless as well.

Creating Symmetry/Invariance via Minimavity. Pivots are most useful when constructing confidence intervals (or more generally, confidence sets) whose performance is invariant to f_θ . We can try to apply the same logic to symmetrize point estimators. This often works when a problem has inherent group symmetries (see Lehmann and Casella, 1998, Chapter 3) but success is no longer guaranteed. This calls for a more general approach. As motivation, consider an asymmetric point estimator. It is asymmetric because it is more accurate for some values of θ than others. To make it more symmetric, we should “redistribute our point estimator’s effort” from those values of θ for which its risk is small to those where its risk is large, thus making its risk as uniform in θ as possible. Just as wealth redistribution helps the poorest in society, this risk redistribution will make the resulting estimator perform better in worst-case scenarios (but as a tradeoff, worse in best-case scenarios).

The above reasoning hints that (approximately) symmetric estimators are those with good performance in worst case scenarios. Therefore, in our search for symmetric estimators, we might focus attention on so-called *minimax* estimators, $\hat{\theta}_{\text{mini}}$, which minimize the average error over $\hat{\Omega}$ under some “least favorable” f_θ specification; see Brown (1994) and Strawderman (2000). However, there are caveats:

1. Minimavity generates invariance through redistribution which can lead to inefficiency. For example, suppose we want to estimate the probability, θ , that a coin lands heads using n flips of the coin. Under squared error loss, the fraction of heads out of n flips (the sample mean) is not a minimax estimator for θ —it fares quite well when the true θ value is near 0 or 1 but poorly when θ is near $\frac{1}{2}$. The actual minimax estimator “redistributes effort” towards $\frac{1}{2}$, giving itself a slight edge over the sample mean for a small interval of θ values near $\frac{1}{2}$ but at a huge cost to accuracy outside this interval (see Lehmann and Casella, 1998, Chapter 5 Example 1.7). In general, *engineered symmetry/invariance*, where none inherently exists, can lead to serious efficiency loss across large regions of the parameter space. In contrast, when a problem contains inherent group symmetries (Lehmann and Casella, 1998, Chapter 3), the redistribution strategy works quite well.
2. As with pivots, the more we match, the harder it becomes to apply the minimax strategy. For example, minimax estimators over the set of complete matches ($D' = D$) do not usually exist because the worst case loss is usually infinity (or the maximum allowable loss) for *any* estimator. Even with partial matching, the price of creating invariance may be too high (see Brown, 1990)—this is simply the relevance robustness tradeoff at play.

3.3 An Extreme Case of Asymmetry

Hypothesis testing problems often display a form of asymmetry so acute that we cannot robustly assess test accuracy even when no matching is involved. That is, we cannot reliably predict whether a hypothesis test will give the right conclusion (at least not without external information).

Example 8 (The Trouble with Testing). Suppose noisy measurements, $\{y_i\}_{i=1}^n$, are taken of the position of a star along some axis, and an additive model, $y_i = \theta + \varepsilon_i$, with ε_i treated as independent standard normal errors, is judged appropriate. Astronomers wish to test the null hypothesis $H_0 : \theta = 0$ against the alternative hypothesis $H_1 : \theta = \theta_1$ where θ_1 is an *unknown, non-zero* value. The canonical z-test looks at the magnitude of the average position, $|\bar{y}_n|$. If the average position is far from 0, say $|\bar{y}_n| > 1.96\sqrt{n}$, the test rejects H_0 . We examine the error rate of this test for various choices of $\hat{\Omega}$.

Setting 1 For all control problems, set $\theta' = 0$ (i.e., we specify \hat{f}_θ as a “point mass” at 0). Our test commits an error if it rejects the null. Under our normal model, $\bar{y}'_n > 1.96\sqrt{n}$, less than 5% of the time when $\theta' = 0$. Thus our test has 5% error rate over $\hat{\Omega}$, its *Type I error* (false positive rate).

Setting 2 For all control problems, set $\theta' = c$ (i.e. we specify \hat{f}_θ as a “point mass” at some non-zero c). Our test now commits an error if it *fails to reject* the null (a Type II error or false negative). The error rate is highly dependent on our choice of the simulation parameter c . As $|c|$ approaches 0, the error rate approaches 95% (cf. Setting 1). For $|c|$ large, the error rate approaches 0%. In practice, we can draw a *power curve* (one minus the error rate plotted as a function of c).

We see that the accuracy of the z-test depends critically on our specification of $\hat{\Omega}$; for arbitrary specifications of \hat{f}_θ , the error rate of our test can be as high as 95% or as low as 0%—an almost vacuous statement.

To understand the source of the asymmetry, note that we chose the $1.96\sqrt{n}$ threshold so that our test would have 5% error rate under Setting 1. That is, hypothesis tests are constructed so that they have guaranteed good performance for select choices of \hat{f}_θ . This “favoritism” for select \hat{f}_θ does not usually exist in interval or point estimation problems and is the reason why the real-life accuracy of a hypothesis test is so hard to predict. The lower our error rate when $\theta = 0$, the higher it must be when θ is near but not equal to 0. To predict the test’s accuracy therefore requires prior information about the magnitude of θ ; without this information, the Type I error and power curve of a test—the two traditional measures of accuracy—do not allow us to predict how well our test will *actually* perform.

Remark 4 *Some of the asymmetry above can be alleviated by testing only for “interesting” alternatives. For example, we may only care about practically different deviations from 0, say $|\theta_1| \geq 1$. With a sample size of 10, restriction of \hat{f}_θ to interesting alternatives bounds the error rate of our test below 11.5%; see Johnson and Rossell (2010) and references therein for more on this approach.*

Except in special cases (see Berger et al., 1994; Berger, 2003), inference for the accuracy of a hypothesis test will not be robust to \hat{f}_θ . These special cases usually require us to treat false rejections on an equal footing with failures to reject, alleviating the asymmetry in the problem. There is hope, however. We are now in an era of massive parallel testing, e.g., millions of genes are tested for association with disease risk. Together, these parallel tests form an “empirical” set of control problems through which we can assess a test’s accuracy (Benjamini and Hochberg, 1995; Efron, 2010). As we will see in Section 4.2, data structures like this open up new ways to achieve relevance and robustness in testing.

4 From the Bayes-Frequentist Dichotomy to the Relevance-Robustness Continuum

4.1 Let’s Compromise: Partial Conditioning

A treatment behaves differently on different patients—this is the motivation for *personalized* medicine. Inspired by this idea, our central theme has been conducting *individualized* inference for a dataset. For example, we can use features of the data, via matching, to assess how well the 5% nominal error rate of a 95% confidence interval applies to the problem at hand. The success of individualization relies on our ability to build controls which resemble our actual problem—this in turn relies on our ability to meaningfully specify \hat{f}_θ and \hat{f}_ε . The greater the level of individualization we strive for, the more stress we put on our model, the more fragile our inference. The degree of individualization should therefore scale with the reliability of our model. This logic underlies the hardest question in statistics: *how individualized* should (can) our inference be?

The difference between (subjective) Bayesian and Frequentist inference hinges on this question. While they are often thought to be two different methodologies, in fact they share the same logic, with the only difference being how they select the relevant subset of control problems. What is usually dubbed as Bayesian statistics, we call full individualization (complete matching); what is dubbed (unconditional) Frequentist statistics, we call no individualization. We purposely eschew the traditional terminology to highlight the fallacy of the dichotomy; Bayes

and Frequentism are two ends of the same spectrum—a spectrum defined in terms of relevance and robustness. The nominal contrast between them—parameters are fixed for Frequentists but random for Bayesians—is a red herring. With this in mind, we note that it appears strange that the two most common answers to the question—how individualized should we be—have also been the *two most extreme answers*. The appropriate tradeoff most likely lies in between, and will be problem-dependent (the principle of individualization!). That is, a compromise between Frequentist and Bayesian inference through partial matching may yield a more satisfactory relevance robustness tradeoff.

Example 9 (Partial Matching). Suppose we have background and outcome information for n individuals, $\{\bar{x}_i, y_i\}_{i=1}^n$, and wish to predict the outcome y_0 for a target individual with background \bar{x}_0 . Hence, $\theta = y_0$ and $D = \{\bar{x}_0, \{\bar{x}_i, y_i\}_{i=1}^n\}$. One possibility is to fit a regression model $y_i = \bar{x}_i^\top \beta + \varepsilon_i$ via least squares, then predict y_0 using $\hat{y}_0 = \bar{x}_0^\top \hat{\beta}$, where $\hat{\beta}$ is the regression estimate. To assess this procedure’s accuracy, we estimate its error, $\Delta = \hat{y}_0 - y_0$, when applied to the target individual.

We simulate controls (D', θ') using the regression model $y'_i = \bar{x}_i^\top \beta' + \varepsilon'_i$, with β' simulated according to a distribution \hat{f}_β and ε'_i as independent standard normal errors. The prediction error for each control problem decomposes into two terms

$$\Delta' = \hat{y}'_0 - y'_0 = \bar{x}_0^\top (\hat{\beta}' - \beta') - \varepsilon'_0. \quad (3)$$

Regardless of which aspects of the data we match on and regardless of our choice for \hat{f}_β , the distribution of the second term ε'_0 is always standard normal. Thus, in deciding the appropriate degree of matching, it suffices to study the first term. Consider now three choices for the set of relevant controls, $\hat{\Omega}_{\text{rele}}$, occupying three positions on the relevance robustness spectrum.

- I. *Complete Matching (Bayes):* $D' = D$. Due to matching, $\hat{\beta}' = \hat{\beta}$ for all controls in $\hat{\Omega}_{\text{rele}}$. The first term in (3) becomes $\bar{x}^\top (\hat{\beta} - \beta')$; its distribution over the set of complete matches (known as the posterior distribution) is sensitive to \hat{f}_β . If we trust \hat{f}_β , however, it is the “optimal” inference (Section 2.2).
- II. *No Matching (Frequentist):* $\hat{\Omega}_{\text{rele}} = \hat{\Omega}$. Without matching, the first term in (3) becomes a pivotal quantity, i.e., its distribution over $\hat{\Omega}$ is independent of \hat{f}_β .

Finally, we strike a compromise by conducting our *error assessment* on “partially” matched controls. For algebraic simplicity, we assume below that individual 1 and the target individual have the same background: $\bar{x}_1 = \bar{x}_0$. Disclaimer: we chose this scheme because the resulting inference, (4), cleanly illustrates how partial matching constitutes a “compromise” inference. This was a pedagogical *not* a practical choice.

- III. *Partial Matching:* D' matches D in all aspects except y'_1 may not equal y_1 . Over this subset of partial matches, it can be shown that the first term in (3) becomes (see Remark 5 below)

$$\bar{x}_0^\top (\hat{\beta}' - \beta') = (1 - h_0) \cdot B' + h_0 \cdot F' \quad (4)$$

where $0 \leq h_0 \leq 1$ acts as a weighting factor. The distribution of B' over $\hat{\Omega}_{\text{rele}}$ turns out to be very close to the posterior distribution formed under Strategy I. If \hat{f}_β is correctly specified, it gives a near optimal characterization of our uncertainty regarding Δ ; as with the Bayes inference above, it is relevant but not robust. The distribution of F' over $\hat{\Omega}_{\text{rele}}$ is independent of \hat{f}_β (cf. Strategy II); it gives a robust but less relevant inference. Partial matching creates a compromise inference by using a weighted average of B' and F' !

Remark 5 The weight h_0 equals $\bar{x}_0^\top (X_n^\top X_n)^{-1} \bar{x}_0$ where X_n is the usual design matrix formed from $\{\bar{x}_i\}_{i=1}^n$; it is often called the leverage. The term B' equals $\bar{x}_0^\top (\hat{\beta} - \beta')$, which coincides with our error term in the complete matching analysis (note that unlike for fully matched controls, $\hat{\beta}' \neq \hat{\beta}$ for partially matched controls). The term F' equals $\bar{x}_0^\top (\hat{\beta}' - \beta') + (\Delta_1 - \Delta'_1)$; its distribution over the set of partial matches is free of \hat{f}_β . Here $\Delta_1 = \hat{y}_1 - y_1$ is the “prediction error” for the first individual (using the least squares estimator) and Δ'_1 is the analogous quantity for control problems.

Partial matching gives us a flexible way to vary the degree of relevance versus robustness in our inference. Complications do exist (see Berger and Wolpert, 1988)—most salient, which features should we match/condition on and which should we ignore? There is no unique best answer. However, this does not imply that we should abandon matching in general or always match on everything—the doctor in her treatment of Mr. Payne has certainly not adopted such radical positions! Even though the *theory* of partial matching is not even partially complete, partial matching can nonetheless serve as a useful tool in *practice* as long as we understand its aim: to assess whether the uncertainty associated with *our* problem should be greater (data are “unlucky”) or less (data are “lucky”) than the uncertainty of the *typical* problem in $\hat{\Omega}$. The goal is to identify features of the data that capture aspects of data-luck that are sufficiently robust to \hat{f}_θ . As long as one keeps this principle in mind, one does not have to be bound by the dogma of Frequentism *or* Bayes. One can be Frequentist, Bayesian, and many other things, including being a Fiducialist (see Section 4.2).

Whereas there is no generic method for deciding the best relevance robustness tradeoff, we can offer a whole host of examples. Figure 7 depicts a Venn diagram—references on the left give examples of inferences that veer too far towards robustness, and those on the right too far towards relevance. The goal is to stand in the middle, in the “football” (a concept developed by Carl Morris and Joe Blitzstein for their inference course at Harvard). Hopefully, the reader can map their own problem in terms of these “reference” points.

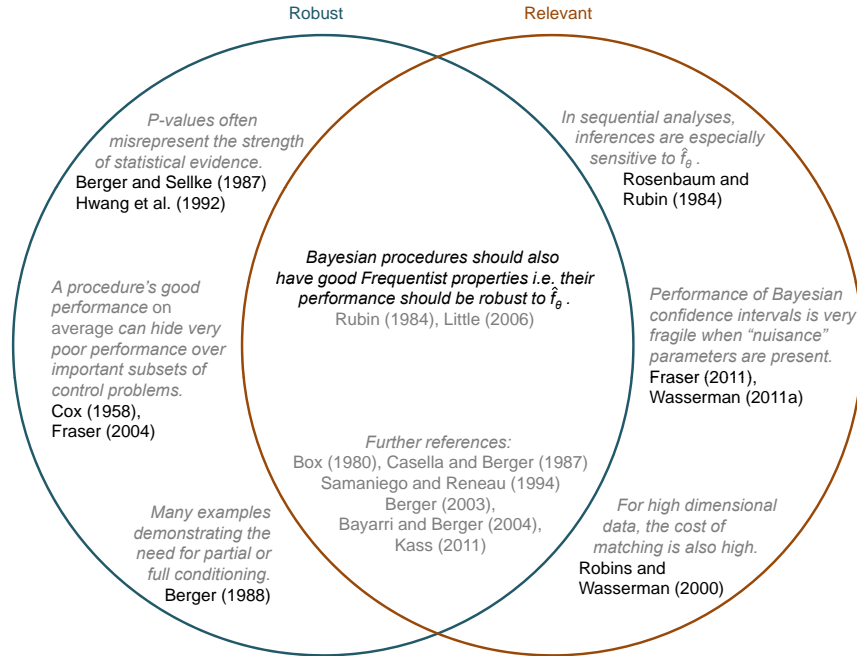


Figure 7: What constitutes a “good” procedure for *our* problem? There are two basic criterion: (a) when we evaluate its performance in a highly individualized manner (i.e., using only controls which closely match our problem), the procedure should have good performance for some reasonable choice of \hat{f}_θ and (b) when we evaluate its performance in a less individualized manner, the procedure should have good performance across a broad range of \hat{f}_θ specifications. Such procedures can be found in the intersection of the Venn diagram. References in the non-intersecting regions discuss the pitfalls of dogmatically skewing in favor of relevance or robustness.

4.2 The Feasibility Diagonal

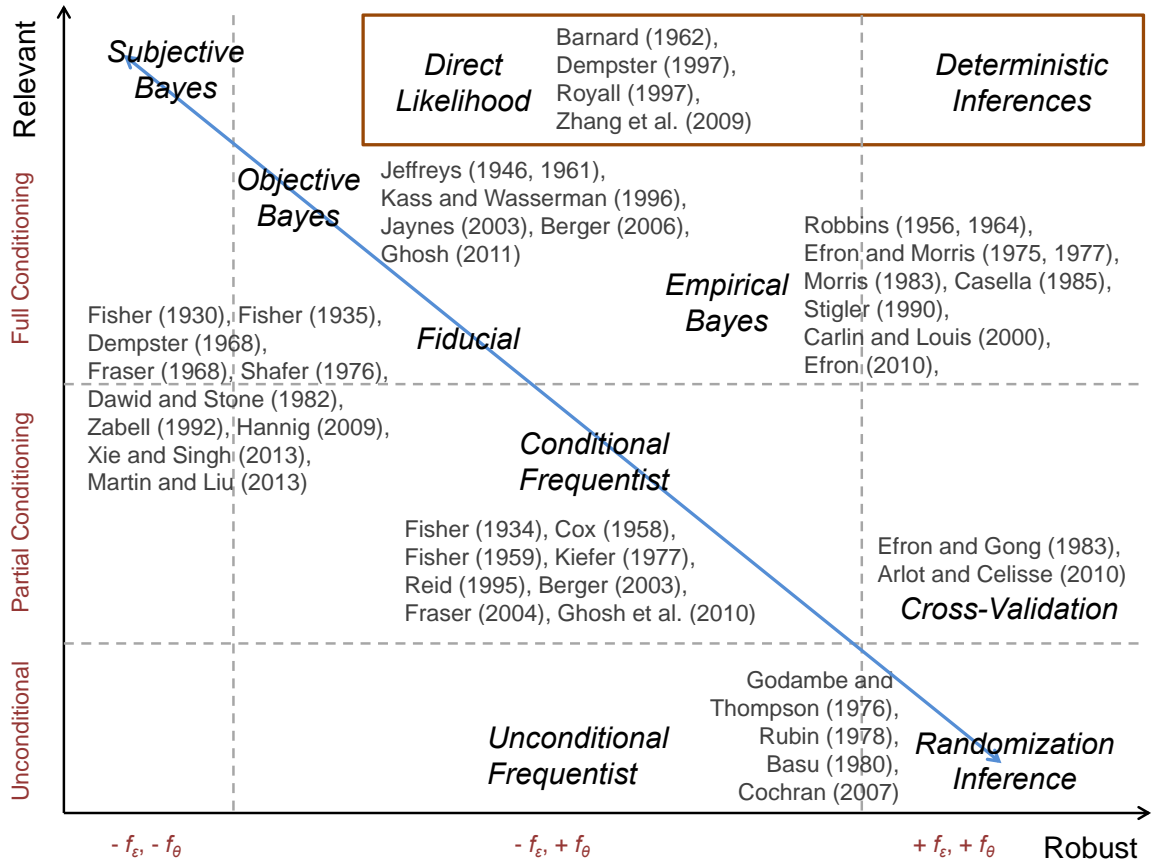


Figure 8: Methods of inference are mapped according to their (robustness, relevance) coordinates. + denotes robustness to modeling of f_θ or f_ε and – denotes a lack of such robustness. The diagonal line is the feasibility diagonal. To move above the diagonal requires additional assumptions or sacrifices. The boxed region in the top-right corner depicts methods which give non-probabilistic inferences.

Our central point is that there is no “right” inference, *only tradeoffs*. Figure 8 maps several common statistical methodologies based on the relevance robustness tradeoff they make. The blue line marks the *feasibility diagonal*—it represents the total budget we can spend on relevance or robustness. At one end of the diagonal lies subjective Bayesian inference, the optimal allocation if we completely trust the information encoded in \hat{f}_θ and \hat{f}_ε . At the other end lies randomization inference (e.g., the Fisher Randomization test; see Basu 1980) whose validity is free of both \hat{f}_θ and \hat{f}_ε . In between lies the objective Bayesian, the Fiducialist, and the conditional Frequentist (A caveat: we have grouped under “Fiducial” several different methodologies whose reasoning process share deep connections with Fisher’s original development). Historically, these three approaches arose out of attempts to create an “objective” inference. Many of these efforts ran into trouble because defining “objectivity” is highly non-trivial (if not impossible)! By situating these three developments along the feasibility diagonal, we may reinterpret the historical quest for “objectivity” as a quest for an optimal relevance robustness tradeoff. Our hope is that the latter perspective offers a more tractable formulation for the problem of inference. Along these lines, one should view objective Bayes, Fiducial and conditional Frequentism as promising initial starting points (having been beta-tested on canonical examples)

in our search for a satisfactory balance between relevance and robustness. But ideally, we will adapt these initial proposals to the special features of the problem at hand. That is, we will follow David Cox’s advice that there are no routine statistical questions but only questionable statistical routines, and discover problem specific sweet spots along the feasibility diagonal.

Moving off the feasibility diagonal. It is trivial to move below the feasibility diagonal (a suboptimal position) by failing to match on important ancillary features of the data (see Example 2). But can we move above the diagonal, to achieve increases in both relevance and robustness? There is of course no free lunch—to move above the diagonal, either a sacrifice must be made (e.g., in terms of interpretability or scope of applicability) or else an assumption must be added. We consider four important examples:

- *Deterministic Inferences*: Statistical inference is overwhelmingly probabilistic because in most situations it is impossible to make non-trivial deterministic statements. However, as Wasserman (2011b) points out, there do exist important applications (for example, see Cesa-Bianchi, 2006) where deterministic control of prediction error can be achieved.
- *Direct Likelihood Inference* (Barnard et al., 1962; Royall, 1997): Likelihood is a way of measuring the degree to which the data supports various hypotheses or parameter values. Importantly, the likelihood depends only on our specification of \hat{f}_ε and not of \hat{f}_θ —thus the allure of drawing inferences directly from the likelihood. Zhang et al. (2009) give an example where this strategy may be appropriate. The sticking point is a lack of interpretability, i.e., “likelihood units” do not have an operational/real life meaning. Thus, while direct likelihood methods (and deterministic inferences) score high on both relevance and robustness, they score low on a third dimension—*utility*— which could not be visualized in our 2-D figure. Remember: no free lunch.
- *Cross Validation* (Efron and Gong, 1983; Arlot and Celisse, 2010): Previously we used the data to build a simulation model for control problems, (D', θ') . Certain applications allow for a more direct approach to creating controls by partitioning the original data as $D = (D', \theta')$. The D' created in this way is known as the “training set”, θ' as the “test set”. Consider Example 9, where we wish to assess the least squares prediction error for a target individual with background \vec{x}_0 . Suppose as before, $\vec{x}_0 = \vec{x}_1$. Hence the first individual is a good proxy for the target individual. This suggests a “leave-one-out” strategy: fit a regression model using only individuals 2 to n (the training set) and compute the error from using the fitted model to predict the first individual’s response (the test set). This error is an estimate of our prediction error for the target individual; it can be a rather noisy estimate because we use a control population of size 1 (the first individual). To reduce our estimate’s variability, we might apply the leave-one-out strategy to *each* of the n individuals in our dataset, i.e., we predict each individual’s response using the remaining $n - 1$ individuals and then average the resulting errors. This effectively gives us a control population of size n . In comparison to Example 9 where we simulated controls using a regression model, the leave-one-out strategy does not require an explicit model—our inference will be more robust. The downside is that we base our accuracy assessment on observed prediction errors for individuals who do not necessarily match the target’s background \vec{x}_0 —our inference will be less relevant.
- *Empirical Bayes* (Robbins, 1956; Efron and Morris, 1977; Casella, 1985): Certain data structures allow for a reliable specification of \hat{f}_θ . Consider Example 6 where the doctor wants to predict your disease status, θ , based on the result of a test, y . In its original form, the problem favors no matching because the doctor may not know the proportion of sick individuals in the population. But now suppose the doctor has an additional source of information: test results $\{y_i\}_{i=1}^n$ from n of his other patients. This allows us to accurately estimate the prevalence, at least when n is large: if the prevalence is low/high, there should be proportionally more negative/positive results among the n patients. Empirical Bayes refers to methods which exploit such forms of indirect information (Efron, 2010) to reliably deliver highly individualized inferences—here, the test results for the *other* n patients provide indirect information about *your* disease status.

Whereas the relevance robustness tradeoff is inescapable, these examples demonstrate how

the particularities of a problem might allow us to recast the exact terms of the tradeoff. It is through such tinkering that innovations in inference take place. Figure 8 depicts the world of inference—relevance and robustness are its longitude and latitude. We were only able to hit the tourist traps on this trip, but we hope the reader, armed with an understanding of relevance and robustness, will return to navigate Figure 8 at their own leisure.

Epilogue: A Rebirth for Individualized Inference?

Standard statistical procedures guarantee good performance *on average* over some set of controls; they are targeted at the “typical” control problem. But every problem is *atypical* in a different way. Just as the doctor must judge whether a clinically successful treatment will continue to be successful given Mr. Payne’s unique background, so the data analyst must gauge how well a procedure’s nominal guarantee, e.g., 95% confidence or 5% Type I error, applies to *this dataset*. Individualized inference aims to deliver uncertainty assessments that are relevant for the problem at hand.

Is Individualization Even Possible? To develop a personalized treatment, we need to conduct personalized experiments, which requires personalized guinea pigs. The heart of statistical inference lies in *designing* relevant control problems. Many mistake the probabilistic statements statisticians make as “natural properties” and statisticians as their “discoverers.” Nothing could be further from the truth. Probability does not simply exist, it is *given existence* through judgements about which control problems are and are not relevant for the actual problem. This point also reveals the risk of individualization—if our judgement is poor, our inferences will be worse off. The quality of our model limits the degree of individualization we can reasonably achieve—to go beyond this limit is to inject misinformation. Divisions in the statistical community center around this question: what is the *appropriate degree*? Frequentists and Bayesians stand at the extremes of no individualization and complete individualization, but the degree of individualization should vary with the problem. The extent to which a doctor personalizes treatment depends on her understanding of the disease and treatment; rather than any fixed relevance robustness tradeoff, the analyst must adapt her inference to the problem environment. Those arguing against individualization may deem such adaptation “ad hoc” and “subjective”—this *mistakes insensitivity for objectivity*.

Why Individualization Matters More Today. A more apt term for Big Data is perhaps complex data. With increasingly complex structures, the ways a dataset can be *atypical* multiply. The statistical models of the past often produced relatively *homogeneous* control problems. The costs of individualization often outweighed its benefits. To accommodate the complex features of modern data, statistical models today inevitably produce ever more *heterogeneous* sets of controls (Efron 2007 provides a striking example in the context of microarray data). The persuasiveness of guarantees on *average performance* deteriorates with increased heterogeneity (see Example 5). Hence the need to deliver *individualized assessments of uncertainty* is more pressing than ever.

Acknowledgements

Deeply insightful comments from Jessica Hwang, Paulo Orenstein, Robin Gong, Andrew Gelman, Iavor Bojinov, Guillaume Basse, and William Meeker improved this paper greatly. The reader has them to thank for any clarities and us to thank for all the inclarities. Keli is indebted to Joe Blitzstein and Carl Morris for their impassioned teaching; they made statistical inference into a beautiful way to reason rather than a series of theorems and formulas. An earlier version of this paper served as the basis for a presentation at the First International Workshop on Bayesian, Fiducial and Frequentist Inference (BFF) held in Shanghai; we appreciate the feedback and encouragements from all the BFFs (Best Friends Forever) in attendance. We also thank the editor Stephen Stigler for his comments and his extraordinary patience as we completed this work. Finally, we thank the NSF and JTF for partial financial support.

References

- Arlot, S. and A. Celisse (2010). A survey of cross-validation procedures for model selection. *Statistics Surveys* 4, 40–79.
- Barnard, G. A., G. M. Jenkins, and C. B. Winsten (1962). Likelihood inference and time series. *Journal of the Royal Statistical Society, A* 125(3), 321–372.
- Basu, D. (1980). Randomization analysis of experimental data: The Fisher randomization test. *Journal of the American Statistical Association* 75(371), 575–582.
- Bayarri, M. J. and J. O. Berger (2000). P values for composite null models. *Journal of the American Statistical Association* 95(452), 1127–1142.
- Bayarri, M. J. and J. O. Berger (2004). The interplay of Bayesian and Frequentist analysis. *Statistical Science* 19(1), 58–80.
- Benjamini, Y. and Y. Hochberg (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society, B* 57(1), 289–300.
- Berger, J. (1985). The Frequentist viewpoint and conditioning. In L. LeCam and R. Olshen (Eds.), *Proceedings of the Berkeley Conference in Honor of Jerzy Neyman and Jack Kiefer*, Volume 1, pp. 15–44.
- Berger, J. O. (1988). An alternative: The estimated confidence approach. In S. S. Gupta and J. O. Berger (Eds.), *Statistical Decision Theory and Related Topics IV*, pp. 85–90.
- Berger, J. O. (2003). Could Fisher, Jeffreys and Neyman have agreed on testing? *Statistical Science* 18(1), 1–32.
- Berger, J. O. (2006). The case for objective Bayesian analysis. *Bayesian Analysis* 1(3), 385–402.
- Berger, J. O., L. D. Brown, and R. L. Wolpert (1994). A unified conditional Frequentist and Bayesian test for fixed and sequential simple hypothesis testing. *The Annals of Statistics* 22(4), 1787–1807.
- Berger, J. O. and T. Sellke (1987). Testing a point null hypothesis: the irreconcilability of P values and evidence. *Journal of the American Statistical Association* 82(397), 112–122.
- Berger, J. O. and R. L. Wolpert (1988). *The Likelihood Principle*. Institute of Mathematical Statistics.
- Box, G. E. (1980). Sampling and Bayes’ inference in scientific modelling and robustness. *Journal of the Royal Statistical Society, A* 143(4), 383–430.
- Breiman, L. (2001). Statistical modeling: The two cultures (with discussion). *Statistical Science* 16(3), 199–231.
- Brown, L. D. (1978). A contribution to Kiefer’s theory of conditional confidence procedures. *The Annals of Statistics* 6(1), 59–71.
- Brown, L. D. (1990). An ancillarity paradox which appears in multiple linear regression (with discussion). *The Annals of Statistics* 18(2), 471–493.
- Brown, L. D. (1994). Minimavity, more or less. In S. S. Gupta and J. O. Berger (Eds.), *Statistical Decision Theory and Related Topics V*, pp. 1–18. Springer.
- Burnham, K. P. and D. R. Anderson (2002). *Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach*. Springer.

- Carlin, B. P. and T. A. Louis (2000). *Bayes and Empirical Bayes Methods for Data Analysis*. CRC Press.
- Casella, G. (1985). An introduction to empirical Bayes data analysis. *The American Statistician* 39(2), 83–87.
- Casella, G. (1992). Conditional inference from confidence sets. In M. Ghosh and P. K. Pathak (Eds.), *Current Issues in Statistical Inference: Essays in Honor of D. Basu*, pp. 1–12. Institute of Mathematical Statistics.
- Casella, G. and R. L. Berger (1987). Reconciling Bayesian and frequentist evidence in the one-sided testing problem. *Journal of the American Statistical Association* 82(397), 106–111.
- Cesa-Bianchi, N. (2006). *Prediction, Learning, and Games*. Cambridge University Press.
- Cochran, W. G. (2007). *Sampling Techniques*. John Wiley & Sons.
- Conti, S., J. P. Gosling, J. E. Oakley, and A. O’Hagan (2009). Gaussian process emulation of dynamic computer codes. *Biometrika* 96(3), 663–676.
- Cox, D. R. (1958). Some problems connected with statistical inference. *Annals of Mathematical Statistics* 29(2), 357–372.
- Cox, D. R. (1982). Randomization and concomitant variables in the design of experiments. In G. Kallianpur, P. R. Krishnaiah, and J. K. Ghosh (Eds.), *Statistics and Probability: Essays in Honor of C. R. Rao*, pp. 197–202. North-Holland.
- Cox, D. R. (2006). *Principles of Statistical Inference*. Cambridge University Press.
- Daubechies, I. (2010). Wavelets and applications. In T. Gowers, J. Barrow-Green, and I. Leader (Eds.), *The Princeton Companion to Mathematics*, pp. 848–862.
- Davison, A. C. and D. V. Hinkley (1997). *Bootstrap Methods and Their Application*. Cambridge University Press.
- Dawid, A. P. and M. Stone (1982). The functional-model basis of fiducial inference. *The Annals of Statistics* 10(4), 1054–1067.
- Dempster, A. P. (1968). A generalization of Bayesian inference. *Journal of the Royal Statistical Society, B* 30(2), 205–247.
- Dempster, A. P. (1997). The direct use of likelihood for significance testing. *Statistics and Computing* 7(4), 247–252.
- Donoho, D. L., I. M. Johnstone, G. Kerkyacharian, and D. Picard (1995). Wavelet shrinkage: Asymptopia? *Journal of the Royal Statistical Society, B* 57(2), 301–369.
- Efron, B. (1986). Why isn’t everyone a Bayesian? *The American Statistician* 40(1), 1–5.
- Efron, B. (1998). R. A. Fisher in the 21st century. *Statistical Science* 13(2), 95–114.
- Efron, B. (2007). Correlation and large-scale simultaneous significance testing. *Journal of the American Statistical Association* 102(477), 93–103.
- Efron, B. (2010). *Large-Scale Inference: Empirical Bayes Methods for Estimation, Testing, and Prediction*. Cambridge University Press.
- Efron, B. and G. Gong (1983). A leisurely look at the bootstrap, the jackknife, and cross-validation. *The American Statistician* 37(1), 36–48.
- Efron, B. and C. Morris (1975). Data analysis using Stein’s estimator and its generalizations. *Journal of the American Statistical Association* 70(350), 311–319.

- Efron, B. and C. Morris (1977). Stein's paradox in statistics. *Scientific American* 236(5), 119–127.
- Efron, B. and R. J. Tibshirani (1994). *An Introduction to the Bootstrap*. CRC Press.
- Evans, M. and H. Moshonov (2006). Checking for prior-data conflict. *Bayesian Analysis* 1(4), 893–914.
- Fisher, R. A. (1925). Theory of statistical estimation. *Mathematical Proceedings of the Cambridge Philosophical Society* 22(5), 700–725.
- Fisher, R. A. (1930). Inverse probability. *Mathematical Proceedings of the Cambridge Philosophical Society* 26(4), 528–535.
- Fisher, R. A. (1934). Two new properties of mathematical likelihood. *Proceedings of the Royal Society of London, A* 144(852), 285–307.
- Fisher, R. A. (1935). The fiducial argument in statistical inference. *Annals of Eugenics* 6(4), 391–398.
- Fisher, R. A. (1959). Mathematical probability in the natural sciences. *Technometrics* 1(1), 21–29.
- Fourdrinier, D. and M. T. Wells (2012). On improved loss estimation for shrinkage estimators. *Statistical Science* 27(1), 61–81.
- Fraser, D. A. S. (1968). *The Structure of Inference*. John Wiley & Sons.
- Fraser, D. A. S. (2004). Ancillaries and conditional inference. *Statistical Science* 19(2), 333–369.
- Fraser, D. A. S. (2011). Is Bayes posterior just quick and dirty confidence? *Statistical Science* 26(3), 299–316.
- Fraser, D. A. S., N. Reid, E. Marras, and G. Yi (2010). Default priors for Bayesian and frequentist inference. *Journal of the Royal Statistical Society, B* 72(5), 631–654.
- Gelman, A., J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin (2013). *Bayesian Data Analysis*. CRC Press.
- Gelman, A. and J. Hill (2006). *Data Analysis Using Regression and Multilevel/Hierarchical Models*. Cambridge University Press.
- Ghosh, M. (2011). Objective priors: An introduction for frequentists. *Statistical Science* 26(2), 187–202.
- Ghosh, M., N. Reid, and D. A. S. Fraser (2010). Ancillary statistics: A review. *Statistica Sinica* 20(4), 1309–1332.
- Godambe, V. P. and M. E. Thompson (1976). Philosophy of survey-sampling practice. In W. Harper and C. Hooker (Eds.), *Foundations of Probability Theory, Statistical Inference, and Statistical Theories of Science*, pp. 103–123. Springer.
- Goutis, C. and G. Casella (1995). Frequentist post-data inference. *International Statistical Review* 63(3), 325–344.
- Hannig, J. (2009). On generalized fiducial inference. *Statistica Sinica* 19(2), 491–544.
- Hansen, M. H. and B. Yu (2001). Model selection and the principle of minimum description length. *Journal of the American Statistical Association* 96(454), 746–774.
- Hwang, J. T., G. Casella, C. Robert, M. T. Wells, and R. H. Farrell (1992). Estimation of accuracy in testing. *The Annals of Statistics* 20(1), 490–509.

- Ioannidis, J. P. A., E. E. Ntzani, T. A. Trikalinos, and D. G. Contopoulos-Ioannidis (2001). Replication validity of genetic association studies. *Nature genetics* 29(3), 306–309.
- Jaynes, E. T. (2003). *Probability Theory: The Logic of Science*. Cambridge University Press.
- Jeffreys, H. (1946). An invariant form for the prior probability in estimation problems. *Proceedings of the Royal Society of London, A* 186(1007), 453–461.
- Jeffreys, H. (1961). *Theory of Probability* (3 ed.). Oxford University Press.
- Johnson, V. E. and D. Rossell (2010). On the use of non-local prior densities in Bayesian hypothesis tests. *Journal of the Royal Statistical Society, B* 72(2), 143–170.
- Kass, R. E. (2011). Statistical inference: The big picture. *Statistical Science* 26(1), 1–9.
- Kass, R. E. and L. Wasserman (1996). The selection of prior distributions by formal rules. *Journal of the American Statistical Association* 91(435), 1343–1370.
- Kennedy, M. C. and A. O’Hagan (2001). Bayesian calibration of computer models. *Journal of the Royal Statistical Society, B* 63(3), 425–464.
- Kiefer, J. (1976). Admissibility of conditional confidence procedures. *The Annals of Statistics* 4(5), 836–865.
- Kiefer, J. (1977). Conditional confidence statements and confidence estimators. *Journal of the American Statistical Association* 72(360a), 789–808.
- Kruskal, W. (1988). Miracles and statistics: The casual assumption of independence. *Journal of the American Statistical Association* 83(404), 929–940.
- Leeb, H. and B. M. Pötscher (2005). Model selection and inference: Facts and fiction. *Econometric Theory* 21(01), 21–59.
- Lehmann, E. L. (1990). Model specification: the views of Fisher and Neyman, and later developments. *Statistical Science* 5(2), 160–168.
- Lehmann, E. L. and G. Casella (1998). *Theory of Point Estimation* (2 ed.). Springer.
- Little, R. J. (2006). Calibrated Bayes: a Bayes/frequentist roadmap. *The American Statistician* 60(3), 213–223.
- Lu, K. L. and J. O. Berger (1989). Estimation of normal means: frequentist estimation of loss. *The Annals of Statistics* 17(2), 890–906.
- Martin, R. and C. Liu (2013). Inferential models: A framework for prior-free posterior probabilistic inference. *Journal of the American Statistical Association* 108(501), 301–313.
- Meng, X.-L. (2014). A trio of inference problems that could win you a Nobel Prize in statistics (if you help fund it). In X. Lin, D. L. Banks, C. Genest, G. Molenberghs, D. W. Scott, and J.-L. Wang (Eds.), *In The Past, Present and Future of Statistical Science*, pp. 535–560. CRC Press.
- Morris, C. N. (1983). Parametric empirical Bayes inference: theory and applications. *Journal of the American Statistical Association* 78(381), 47–55.
- Neyman, J. (1934). On the two different aspects of the representative method: the method of stratified sampling and the method of purposive selection. *Journal of the Royal Statistical Society* 97(4), 558–625.
- Reid, N. (1995). The roles of conditioning in inference. *Statistical Science* 10(2), 138–157.

- Robbins, H. (1956). An empirical Bayes approach to statistics. In *Proceedings of the Third Berkeley Symposium on Mathematics, Statistics and Probability*, Volume 1, pp. 157–163. University of California Press, Berkeley.
- Robbins, H. (1964). The empirical Bayes approach to statistical decision problems. *The Annals of Mathematical Statistics* 35(1), 1–20.
- Robins, J. and L. Wasserman (2000). Conditioning, likelihood, and coherence: A review of some foundational concepts. *Journal of the American Statistical Association* 95(452), 1340–1346.
- Robinson, G. K. (1979a). Conditional properties of statistical procedures. *The Annals of Statistics* 7(4), 742–755.
- Robinson, G. K. (1979b). Conditional properties of statistical procedures for location and scale parameters. *The Annals of Statistics* 7(4), 756–771.
- Rosenbaum, P. R. (1984). Conditional permutation tests and the propensity score in observational studies. *Journal of the American Statistical Association* 79(387), 565–574.
- Rosenbaum, P. R. and D. B. Rubin (1984). Sensitivity of Bayes inference with data-dependent stopping rules. *The American Statistician* 38(2), 106–109.
- Royall, R. (1997). *Statistical Evidence: A Likelihood Paradigm*. CRC Press.
- Rubin, D. B. (1978). Bayesian inference for causal effects: The role of randomization. *The Annals of Statistics* 6(1), 34–58.
- Rubin, D. B. (1984). Bayesianly justifiable and relevant frequency calculations for the applied statistician. *The Annals of Statistics* 12(4), 1151–1172.
- Rukhin, A. L. (1988). Estimated loss and admissible loss estimators. In S. S. Gupta and J. O. Berger (Eds.), *Statistical Decision Theory and Related Topics IV*, pp. 409–418.
- Samaniego, F. J. and D. M. Reneau (1994). Toward a reconciliation of the Bayesian and frequentist approaches to point estimation. *Journal of the American Statistical Association* 89(427), 947–957.
- Shafer, G. (1976). *A mathematical theory of evidence*. Princeton University Press.
- Stigler, S. M. (1990). The 1988 Neyman memorial lecture: a Galtonian perspective on shrinkage estimators. *Statistical Science* 5(1), 147–155.
- Strawderman, W. E. (2000). Minimavity. *Journal of the American Statistical Association* 95(452), 1364–1368.
- Sundberg, R. (2003). Conditional statistical inference and quantification of relevance. *Journal of the Royal Statistical Society, B* 65(1), 299–315.
- Wasserman, L. (2011a). Frasian inference. *Statistical Science* 26(3), 322–325.
- Wasserman, L. (2011b). Low assumptions, high dimensions. *Rationality, Markets and Morals* 2(49), 201–209.
- Xie, M.-g. and K. Singh (2013). Confidence distribution, the frequentist distribution estimator of a parameter: A review. *International Statistical Review* 81(1), 3–39.
- Zabell, S. L. (1992). R. A. Fisher and the fiducial argument. *Statistical Science* 7(3), 369–387.
- Zhang, J. L., D. B. Rubin, and F. Mealli (2009). Likelihood-based analysis of causal effects of job-training programs using principal stratification. *Journal of the American Statistical Association* 104(485), 166–176.