

2009

Accounting for Differences among Patients in the FDA Approval Process

Mark Van Der Laan

Anup Malani

Oliver Van Der Bembom

Follow this and additional works at: [https://chicagounbound.uchicago.edu/
public_law_and_legal_theory](https://chicagounbound.uchicago.edu/public_law_and_legal_theory)

 Part of the [Law Commons](#)

Chicago Unbound includes both works in progress and final versions of articles. Please be aware that a more recent version of this article may be available on Chicago Unbound, SSRN or elsewhere.

Recommended Citation

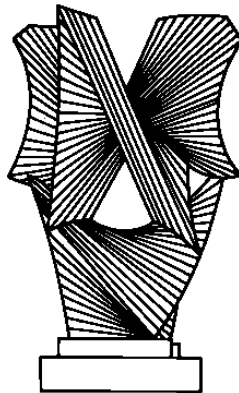
Mark van der Laan, Anup Malani & Oliver Van Der Bembom, "Accounting for Differences among Patients in the FDA Approval Process" (University of Chicago Public Law & Legal Theory Working Paper No. 281, 2009).

This Working Paper is brought to you for free and open access by the Working Papers at Chicago Unbound. It has been accepted for inclusion in Public Law and Legal Theory Working Papers by an authorized administrator of Chicago Unbound. For more information, please contact unbound@law.uchicago.edu.

CHICAGO

JOHN M. OLIN LAW & ECONOMICS WORKING PAPER NO. 488
(2D SERIES)

PUBLIC LAW AND LEGAL THEORY WORKING PAPER NO. 281



ACCOUNTING FOR DIFFERENCES AMONG PATIENTS IN THE FDA APPROVAL PROCESS

Anup Malani, Oliver Bembom and Mark van der Laan

THE LAW SCHOOL
THE UNIVERSITY OF CHICAGO

October 2009

This paper can be downloaded without charge at the John M. Olin Program in Law and Economics Working Paper Series: <http://www.law.uchicago.edu/Lawecon/index.html> and at the Public Law and Legal Theory Working Paper Series: <http://www.law.uchicago.edu/academics/publiclaw/index.html> and The Social Science Research Network Electronic Paper Collection.

Accounting for Differences among Patients in the FDA Approval Process

Anup Malani, Oliver Bembom and Mark van der Laan¹

Abstract. The FDA employs an average-patient standard when reviewing drugs: it approves a drug only if the average patient (in clinical trials) does better on the drug than on control. It is common, however, for different patients to respond differently to a drug. Therefore, the average-patient standard can result in approval of a drug with significant negative effects for certain patient subgroups (false positives) and disapproval of drugs with significant positive effects for other patient subgroups (false negatives). Drug companies have a financial incentive to avoid false negatives. After their clinical trials reveal that their drug does not benefit the average patient, they conduct what is called *post hoc* subgroup analysis to highlight patients that benefit from the drug. The FDA rejects such analysis due to the risk of spurious results. With enough data dredging, a drug company can always find some patients that benefit from their drug.

This paper asks whether there workable compromise between the FDA and drug companies. Specifically, we seek a drug approval process that can use *post hoc* subgroup analysis to eliminate false negatives but does not risk opportunistic behavior and spurious correlation. We recommend that the FDA or some other independent agent conduct subgroup analysis to identify patient subgroups that may benefit from a drug. Moreover, we suggest a number of statistical algorithms that operate as veil of ignorance rules to ensure that the independent agent is not indirectly captured by drug companies. We illustrate our proposal by applying it to the results of a recent clinical trial of a cancer drug (motexafin gadolinium) that was recently rejected by the FDA.

¹ University of Chicago Law School; Department of Statistics, University of California-Berkeley; and Department of Statistics, University of California-Berkeley, respectively. We thank Richard Miller for access to the data from the Phase III motexafin gadolinium (MgD) trials. We thank Glenn Cohen, Adam Cox, Jack Calfee, Richard Epstein, Jake Gersen, Jonathan Masur, Richard Miller, Lars Noah, Ben Roin, participants at the AEI conference, "Oncology Drug Development: Rethinking FDA Oversight," March 13-14, 2008, and workshop participants at the University of Chicago Law School for helpful comments. Malani thanks the Milton and Miriam Handler Foundation for financial support. In the interest of full disclosure, two of the authors, Bembom and van der Laan work for Target Analytics, Inc., which was hired by an investor in Pharmacyclics to prepare a statistical analysis of the data from the MgD trials. The data analysis for this particular paper was conducted after Pharmacyclics' NDA for MgD was rejected by the FDA in Dec. 2007. B. Jungbauer, Pharmacyclics' Xcytrin Gets FDA "Not Approvable" For NSCLC Patients With Brain Metastases, The Pink Sheet (Dec. 28, 2007).

Introduction

Different people respond differently to drugs. Statisticians refer to this as heterogeneity in treatment response. For example, Arthrotec (Pfizer) is an effective treatment for osteoarthritis for patients who develop ulcers when certain common pain medications.² But a key ingredient in Arthrotec – misoprostol – is documented to induce labor and used for medical elective abortion.³ Therefore, while Arthrotec is generally effective for pain relief, it is contraindicated⁴ for pregnant women because of its abortifacient effects.

Or consider the case of motexafin gadolinium (MGd), a treatment for lung cancer patients experiencing mental illness because the cancer has spread to – or metastasized in – their brain. In a Phase III trial, the average patient did not experience a statistically significant benefit from the drug.⁵ This failure is driven by poor results among patients who had already taken chemotherapy for their cancer.⁶ Among the patients who were newly diagnosed and had not yet received any chemotherapy, however, the drug appears to more than double the median time before the onset of dementia.⁷

Under current regulations, however, the U.S. Food and Drug Agency (FDA) typically considers only the effect of a drug in the average patient when deciding whether to approve the drug. Given heterogeneity in treatment response, this approach can result in the approval of drugs with significant negative effects for identifiable subgroups of patient (false positives) and in the non-approval of drugs with significant positive effects for identifiable subgroups (false negatives).

The FDA is not entirely deaf to these concerns. To address false positives, the FDA may perform (or require the drug sponsor⁸ to perform) so-called *post hoc* subgroup analysis to identify subgroups that do not benefit or that suffer severe side effects from an otherwise

² Arthrotec is diclofenac sodium plus misoprostol. Osteoarthritis is pain and inflammation caused by the breakdown of cartilage in a patient's joints. The pain medications that cause ulcers are non-steroidal anti-inflammatory drugs (NSAIDs).

³ See Goldberg et al., Misoprostol and Pregnancy, 34(1) New Eng. J. Med. 38 (Jan. 4, 2001).

⁴ Arthrotec product insert at 1, available at FDA, Label and Approval History – Arthrotec, Labeling revision, Aug. 24, 2007 (<http://www.fda.gov/cder/foi/label/2007/020607s0101bl.pdf>) (checked Mar. 7, 2008).

⁵ Statistical significance is conventionally defined as having a p-value less than or equal to 0.05.

⁶ Their tumors tended to be more resilient as they had already survived previous treatment.

⁷ This result is statistically significant, with a p-value less than 0.002. See *infra* Section 4a.

⁸ Any organization conducting clinical trials in order to obtain marketing approval for a drug is called a drug sponsor. A drug company may be a sponsor, but not all sponsors must be drug companies. For example, a non-profit organization such as a University may seek marketing approval for a drug.

approved drug. Subgroup analysis is statistical evaluation of the effect of a drug on one or more subgroups of the subjects in a clinical trial. *Post hoc* subgroup analysis, in particular, is the estimation of treatment effects in subgroups that were not specified prior to the start of the trial. Instead the statistical analysis is performed on subgroups *identified from the data after the trial had begun or was completed*. The FDA may use the results from this *post hoc* subgroup analysis to justify a label that indicates that the drug is inappropriate for the subgroups it identifies.

To address false negatives, the FDA does two things. First, it allows a drug sponsor to specify, before a trial, one or more subsets of the targeted patient population for which it plans to undertake subgroup analysis. In many cases, however, the sponsor may not have enough information prior to trials to identify subgroups that may be especially sensitive to treatment. Even if it did, the FDA would require the sponsor to increase sample size so that its study gathers sufficient statistical information – or “power” – to accurately estimate treatment effects in those subgroups.⁹ This additional cost may outstrip the financial resources available to many sponsors. Second, the sponsor may conduct *post hoc* subgroup analysis following a trial that is not powered for that analysis and ask the FDA for permission to conduct a follow-up trial on a possibly sensitive subgroup identified in the *post hoc* analysis. But this approach adds the expense of a full additional trial to the cost of clinical testing.

Although drug sponsors complain that the FDA’s position imposes too high a cost on exploiting heterogeneity in treatment response, the agency’s approach has a rational foundation. *Post hoc* subgroup analysis increases the risk of approving drugs that have no net beneficial effect. The more subgroups the sponsor analyzes, the more likely it is to find one that appears to benefit even if in fact there exists no subgroup that benefits. We call this problem the risk of spurious correlation from “multiple testing”. The sponsor has a financial incentive to ignore – and perhaps increase – this risk. After all, they have already invested millions in lab research and clinical testing; if the drug is not approved they will obtain no return on this investment. If

⁹ The basic formula for how large the sample *in each subgroup* must be to ascertain a treatment effect is $n = 2\sigma^2(z_{\alpha/2} + z_{\beta})^2/d^2$. In this formula, d is how sensitive the researcher wants the estimate to be. In other words, the formula gives the sample size required to identify treatment effects that are at least as large as d . The formula also depends on σ , the variance of the treatment effect from the drug. The larger the variance, the more the noise in the data and thus the larger the sample size required to identify a treatment effect as small as d . Usually the researcher estimates σ from previous studies of the drug or the disease. The crucial statistical parameters are $z_{\alpha/2}$, which determines defines the confidence level of the analysis, and z_{β} , which determines the power of the analysis. The higher is the confidence level of the analysis, the lower is the chance of a false positive. A confidence level of 95% (i.e., $\alpha = 0.05$) means that the probability that a significant result is false is just 5%. The critical value for this level of confidence in a two-sided test with a normal distribution is $z_{\alpha/2} = 1.96$. The higher is the power of an analysis, the lower is the chance of a false negative. A power level of 80% ($\beta = 0.8$) means that the probability that a drug with a positive treatment effect is mistakenly reported as having an insignificant treatment effect is 20%. The critical value for this level of power is $z_{\beta} = 0.842$. As is apparent from the formula, greater higher confidence levels and power require larger sample size.

the FDA knew the number of subgroups the sponsor sampled in search of a positive response the FDA could limit the risk of spurious correlation by employing statistical corrections for multiple testing, such as raising the p-value¹⁰ required to demonstrate statistical significance.¹¹ Unfortunately, the FDA will rarely be able to verify this number.

This paper asks whether there is a better way. In particular, we explore whether there is a process that allows approval of a drug based on *post hoc* subgroup analysis without the risks of opportunistic behavior and spurious correlation or the cost of unnecessary additional trials. Based on our current statistical understanding, we offer two compromises to achieve this result.

The FDA's current position is that whenever possible sponsors specify the patient subgroups they plan to study *before conducting a trial*. However, there will be cases where treatment-sensitive subgroups are not known before the trial. Our first proposal is to identify these subgroups using an "adaptive group sequential design" trial. In standard trials, patients are typically assigned to the treatment or control group according to a pre-set randomization scheme and remain in their assigned groups. In an adaptive design, new patients can be randomized to treatment or control based on the performance of patients previously enrolled in the trial. The goal is to identify treatment-sensitive subgroups based on data from early enrollees in order to power up analysis of those subgroups among later enrollees. Adaptive trials may require a larger sample size than standard trials, but they do not require as many total subjects as performing a subsequent trial based on *post hoc* analysis following an initial, failed trial. Moreover, the process is faster because only one trial is conducted.

Our second proposal is that, in certain situations, it may be permissible to approve a drug on the basis of *post hoc* analysis if that analysis were done in a manner that eliminated the incentive for data dredging. For example, the analysis could be conducted by an independent consultant rather than the sponsor. To address concerns about whether the consultant is truly independent of the sponsor, we suggest some statistical algorithms that help the consultant identify patient subgroups which the FDA can trust are not the product of data dredging. These algorithms may be thought of as veil of ignorance rules¹² because they blind the consultant to information so as to ensure the consultant is not indirectly captured by the drug sponsor. Our most promising algorithm is what we call "split-sample analysis." This reform would give the consultant access only to a random subsample of the trial data as selected by the FDA. The consultant would be asked to identify subgroups based on *post hoc* subgroup analysis on this

¹⁰ In this context, the p-value of a statistical estimate is one minus the probability that the estimate is different from zero, which signifies no treatment effect. In other words, if $p < 0.05$, then it can be said that we are more than 95% sure that the estimated treatment effect is different than zero.

¹¹ The corrections are explained at greater length in the text accompanying note 38.

¹² For a discussion of this class of rules, see Adrian Vermeule, *Veil of Ignorance Rules in Constitutional Law*, 111 *Yale L. J.* 399 (2001).

“exploratory” sample. The sponsor would then be permitted to seek drug approval for a subgroup without conducting an additional clinical trial under two conditions. First, it can only obtain approval for a subgroup reported by outside consultant. Second, it can only obtain approval for a subgroup if the drug was significantly effective and safe in that subgroup in the remainder of the trial sample, which we call the “confirmatory” sample. Importantly, significance will be judged according to a higher standard to account for the risk of spurious correlation due to multiple testing bias. This protects against false positives in two ways. Because the outside consultant does not have access to confirmatory sample, it cannot help the sponsor out by choosing subgroups that respond positively only in the confirmatory sample. Nor can the consultant help the sponsor by reporting a large number of subgroups based on the exploratory sample because the sponsor pays a multiple-testing penalty in the confirmatory subsample for each additional subgroup that the consultant reports.

To illustrate the potential benefit of increased flexibility in the drug approval process and to demonstrate how our statistical algorithms could address the data manipulation problem, we conduct a case study of Pharmacyclics’ MGd, a drug for the treatment of non-small cell lung cancer (NSCLC) patients with brain metastases. The company was unable to demonstrate efficacy in its confirmatory Phase III trial.¹³ Their own *post hoc* subgroup analysis suggested that the insignificant finding was due to enrollment of subjects whose brain tumors had already proven resistant to radiation therapy. These subjects were unlikely to respond to any treatment, including MGd. Excluding these subjects, the sponsor showed MGd had a statistically significant treatment effect.¹⁴ However, the FDA rejected Pharmacyclics’ application to market the drug.¹⁵ Mimicking the role of outside consultants, we apply our proposed statistical algorithms to correct for bias due to any opportunistic behavior by the sponsor. We find that one of our algorithms would have surely validated Pharmacyclics’ claims and the other – split-sample analysis – would have validated them with probability 0.11. That is, there is at least a one-in-ten chance that truly independent analysis would have led to the approval of MGd.

The remainder of the paper has the following structure. Section 1 explains how the FDA currently handles heterogeneous treatment response. Section 2 discusses the problem of spurious correlation and opportunism associated with *post hoc* subgroup analysis. Section 3 presents trial designs and institutional arrangements coupled with statistical analyses that allow identification of subgroups for which the drug is effective in a manner that limits false positives and the cost of trials. Finally, Section 4 examines data from the final Phase III MGd trial to illustrate the new statistical methods we propose. We acknowledge that this article is not the right venue to present

¹³ See *infra* text accompanying notes 64 and 66.

¹⁴ *Id.* See also Personal communication with then-CEO Richard Miller, Mar. 14, 2008.

¹⁵ See B. Jungbauer, Pharmacyclics’ Xcytrin Gets FDA “Not Approvable” For NSCLC Patients With Brain Metastases, *The Pink Sheet* (Dec. 28, 2007).

formal statistical derivations and discuss in the future statistical research required to deal with the statistical question raised by our proposals in the conclusion.

1. FDA policy on heterogeneity in treatment response

The Food, Drug and Cosmetics Act requires that the FDA verify that a drug is safe and effective before it is approved for marketing as a therapeutic.¹⁶ FDA regulations typically require that a company applying for marketing approval conduct two Phase III trials to demonstrate efficacy and relatively tolerable side effects.¹⁷ Although regulations do not spell out exactly the evidentiary standard to which the FDA holds a new drug, the FDA has issued a guidance that emphasizes a drug should be evaluated based on all the patients that enroll in a trial¹⁸ and that requires the rate of false positive findings be set to 5 percent.¹⁹ The implication – borne out by practice – is that the FDA judges the efficacy of a drug by the difference between average outcomes in the treatment and control arms of a trial.

The FDA understands that there is heterogeneity of treatment effects across patient subgroups.²⁰ But it only accommodates this heterogeneity in a limited way. First, it encourages sponsors to specify prior to conducting a trial – or in statistic parlance, specify *a priori* – the subgroups they plan to analyze.²¹ When this is done, sponsors must account for the risk of spurious correlation due to multiple testing bias when setting their initial sample size and when analyzing data from a trial.²² The FDA guidance does not explicitly state that significant

¹⁶ See 21 U.S.C. § 355(d).

¹⁷ Peter Barton Hutt and Richard A. Merrill, FOOD AND DRUG LAW: CASES AND MATERIALS 527 n. 2 (2d ed. 1991). For a more detailed analysis of the two-trial requirement, see Jennifer Kulynych, Will FDA Relinquish the “Gold Standard” for New Drug Approval? Redefining “Substantial Evidence in the FDA Modernization Act of 1997, 54 Food & Drug L. J. 127, 129-130 (1999) (explaining that a second trial is due to the scientific requirement of replication, that the requirement is occasionally waived, and that biologics are less likely to face this requirement).

¹⁸ This is called the intent-to-treat population. It includes even the subjects that drop out. See Food and Drug Agency, International Conference on Harmonisation; Guidance on Statistical Principles for Clinical Trials, Availability, 63(179) Fed. Reg. 49583, 49593 (Sept. 16, 1998) (henceforth “Statistical Guidance” (§5.2.1 Full Analysis Set). The alternative is to evaluate the drug on the per-protocol population. In this case the drug is evaluated only among subjects who enroll, do not drop out, and follow all trial procedures.

¹⁹ Id. at 49291 (§3.5 Sample Size).

²⁰ Id. at 49589 (§3.2 Multicenter Trials).

²¹ Id. at 49595 (§5.7 Subgroups, Interactions, and Covariates).

²² Id. at 49587 (§2.2.5 Multiple Primary Variables). The agency recognizes that the corrections are less severe where subgroups overlap and therefore produce correlated test statistics. If two subgroups are mutually exclusive, the outcomes across groups are statistically independent. The appropriate adjustment for multiple testing bias is the Bonferroni adjustment, which is described in Section 2. If the two subgroups overlap in part, the outcomes of the

treatment effects among *a priori* specified subgroups can be the basis for drug approval, but it does not rule it out.

Second, the FDA acknowledges that it will not always be possible to identify subgroups *a priori* and that exploratory analysis of trial data may be required to identify subgroups.²³ The FDA's approach to subsequent so-called *post hoc* subgroup analysis depends on the average treatment effect in the full trial population and whether the outcome at issue concerns efficacy or safety.

If the sponsor cannot demonstrate that the average treatment effect is sufficient that the drug will be approved for the full trial population, *post hoc* subgroup analysis by itself cannot be used to obtain approval for a subgroup.²⁴ The FDA has not approved a single drug solely on the basis of *post hoc* subgroup analysis.²⁵ The FDA does permit a sponsor to use *post hoc* subgroup analysis identify a subgroup and then to conduct a subsequent trial on that subgroup to confirm the findings of the subgroup analysis. But another trial can be very costly.²⁶ And there is no indication that the FDA allows sponsors to combine or "pool" the data from the subgroup in the initial trial with data on the subgroup from the subsequent confirmatory trial to establish a significant positive result for the subgroup across the two trials.²⁷ Such a combination would mitigate the sample size requirement – and hence costs – for the confirmatory trial.

Even if the sponsor is able to demonstrate that its drug is on average safe and effective for the full trial population in the initial trial, the FDA may require the sponsor to demonstrate the drug is effective and safe for subgroups defined by the agency in order to validate the results for the full trial population. FDA guidelines identify subgroups defined by centers in multicenter trials as one such check on the consistency of the trial's main results,²⁸ but clinically and

groups will be positively correlated. In that case the failure of a test on one group will contribute to the failure of a test for the overlapping second group. This reduces the risk of spurious results from statistical tests on the second group. Therefore, something weaker than the Bonferroni adjustment is required to correct for multiple testing bias.

²³ Id at 49595 (§5.7).

²⁴ Id. at 49595 (§5.7). See also John Powers et al., FDA Evaluation of Antimicrobials: Subgroup Analysis, Letter to Editor, 126(6) Chest 2298 (June 2005).

²⁵ Aldo P. Maggioni, et al., FDA and CPMP Rulings on Subgroup Analyses, 107 Cardiology 97, 99 (2007).

²⁶ It is difficult to estimate from published data the cost of individual trials but it is possible to estimate the cost of different phases of clinical testing. The out-of-pocket plus opportunity costs of phase I are \$45.7 million, phase II are \$65.1 million, and phase III are \$205.5 million. See DiMasi, et al., supra note **Error! Bookmark not defined.**, at 162 (table 1) and 165 (table 3). If phase III involves just two trials, then the cost of per phase III trial is over \$100 million.

²⁷ We examine this approach in infra note 43.

²⁸ FDA, Statistical Guidance, supra note 18, at 49589 (§3.2)

biologically defined subgroups have also been suggested.²⁹ If certain subgroups do not show efficacy or show side effects, the FDA may require that the drug label state it is indicated only for the subpopulations where it has been demonstrated both effective and safe. For example, following the Val-HeFT trial of 160 mg valsartan for patients with heart failure, the FDA only approved the drug for patients who are intolerant to ACE inhibitors. The reason was that in the full trial population the drug was superior to placebo only with respect to only one (combined mortality and morbidity) of the two outcomes studied.³⁰ (The other outcome was mortality alone.) However, the drug was superior on both outcomes in the non-ACE subgroup.³¹ The FDA may also require the sponsor for a drug with an uneven or uncertain safety profile to conduct Phase IV post-approval trials. If the Phase IV trials reveal dangerous side effects, FDA has the ability to alter a drug's labeling to reflect those risks or to yank a drug from the market.

In short, the FDA takes a conservative and asymmetric approach with respect to subgroup analysis.³² If subgroups are identified prior to trial, they may positively influence approval so long as the sponsor powers the study to address multiple testing concerns. If subgroups cannot be identified prior to trial, *post hoc* subgroup analysis can only be used to negatively influence approval or to justify new, costly trials.³³

It is worth noting that the FDA's approach – judging drugs largely on the basis of average treatment effects – implicitly assumes that doctors are very bad at matching the right patient subgroups to drugs.³⁴ To understand why, observe that the average treatment effect of a drug y_1 (relative to control y_0) is equal to its positive treatment effect among patient subgroups that benefit from the drug (i.e., $y_1 > y_0$) plus its *negative* treatment effects among those who do not (i.e., $y_1 < y_0$):

$$E(y_1 - y_0) = pE(y_1 - y_0 | y_1 > y_0) + (1 - p)E(y_1 - y_0 | y_1 < y_0)$$

where $p = \Pr(y_1 > y_0)$ is the fraction of people who benefit from the drug. However, the drug only harms patients among whom it is contraindicated if doctors give those patients the drug.

²⁹ See Powers et al., *supra* note 24, at 2298.

³⁰ The primary endpoint of a trial is the outcome on which the treatment is judged.

³¹ See Maggioni et al., *supra* note 25, at 99.

³² The European Union's Committee for Proprietary Medicinal Products (CPMP) has a similar policy. *Id.* at 97.

³³ See Richard Wunderink et al., FDA Evaluation of Antimicrobials: Subgroup Analysis, Letter to Editor, 126(6) Chest 2300 (June 2005) (highlighting asymmetric implications of *post hoc* subgroup analysis for FDA approval).

³⁴ See Anup Malani and Feifang Hu, The option value of new therapeutics 14, unpublished manuscript (2004). The bad matching may be because there is no way to determine which patients benefit and which do not, because doctors cannot or do not distinguish between patients that benefit and those that do not, or because doctors give the drug to patients that they know will not benefit from it.

The problem is that the FDA's average-effects rule mathematically assumes this occurs, i.e., that doctors give the drug to every patient even if it harms those patients. If the FDA had more faith in doctors, it would instead estimate the value of a drug solely by its positive effects among the subgroup that benefits from the drug.³⁵ The FDA would not have to worry about harming patients for whom the drug is inappropriate because doctors would not give these patients the drug.

2. Cost and benefits from *post hoc* subgroup analysis

The FDA's conservative position on *post hoc* subgroup analysis is based on concerns that multiple testing creates a risk of spurious correlation, which can result in false positive drug approvals.³⁶ To illustrate, consider the following Statistics 101-type hypothetical. Suppose there is a population that can be divided into 10 mutually exclusive subgroups. For example, if the trial population ranges in age uniformly from 20 to 70, we can divide the group into 10 equally sized five-year age bins: 20-24, 25-30, etc. Suppose also that there is a drug which has no effect on either the full population or on any subgroup, though there is some random variation in observed outcomes either due to the drug or natural progression. For example, we might assume that the treatment effect for each subgroup is a normally distributed random variable with mean equal to zero and variance equal to one. The probability that the drug will be proven effective on the full population with a confidence level of 95% is just 5%.

But if the sponsor seeking approval for our hypothetical drug is permitted to separately test the drug against each of the 10 subgroups, the probability that he will be able to demonstrate efficacy for at least one subgroup is $0.4 (= 1 - 0.95^{10})$. This obviously raises the risk of false positives, that is, the possibility of approving the drug even though it has not been demonstrated effective at the 95% confidence level.

If the FDA knows that the sponsor will test the drug against 10 subgroups, it can implement a multiple testing correction or penalty to eliminate spurious results. For example, if the outcomes in the 10 subgroups are known to be uncorrelated, then it can change the threshold p-value required for approval from $p = 0.05$ to $p = 0.05 / (\text{number of tests})$ or 0.005 .³⁷ This correction is known as the Bonferroni adjustment. It ensures the probability of observing even one subgroup with significant treatment effects is back to 5%. The proper p-value adjustment in

³⁵ This value is larger than the average effect of the drug in clinical trials. Because $E(y_1 - y_0 | y_1 < y_0) < 0$, $pE(y_1 - y_0 | y_1 > y_0) > pE(y_1 - y_0 | y_1 > y_0) + (1 - p)E(y_1 - y_0 | y_1 < y_0)$.

³⁶ See Kulynych, *supra* note 17, at 141, John A. Lewis, *Statistical Issues in the Regulation of Medicine*, 14 *Stat. in Med.* 127, 132 (1995).

³⁷ We are ignoring the 11th test on the full population to keep the numbers simple and because the full population result is positively correlated with each subgroup result.

the cases where outcomes in the subgroups are correlated is different, but this correlation and the adjustment may be derived from the data.³⁸

The problem becomes intractable, however, if the FDA does not know the number of subgroups against which the sponsor will test its drug. In that case the FDA cannot impose a proper multiple testing penalty. This problem becomes exponentially more severe the larger the number of possible subgroups. Suppose the sponsor can also divide the adult population by gender (male/female) and by ethnicity (white/black/Hispanic/other) and the drug still has no treatment effect for any subgroup. Now the available subgroups has jumped from 10 based solely on age to 80 ($= 10 \times 2 \times 4$) based on a combination of age, gender and ethnicity. If the sponsor can cherry-pick a subgroup in which to demonstrate efficacy, the probability it will be find able to find at least one with a p-value less than 0.05 is 0.98 ($= 1 - 0.95^{80}$)! The sponsor has a strong financial incentive to cherry-pick in this manner because the alternative may be not to obtain any return on its investment in the drug. We call this the problem of opportunistic behavior by sponsors.³⁹

The FDA's response to this risk is to base its drug approval decision solely on average treatment effects for the full trial population. This bars any approval based on *post hoc* subgroup analysis without further clinical trials. But this response swings the pendulum of error too far in the opposite direction. Instead of approving some drugs with no treatment effect (false positives), the FDA's conservative policy rejects some drugs with positive treatment effects for some subgroups (false negatives) or increases the costs of drugs if the sponsor conducts a follow-on trial to confirm the results of *post hoc* subgroup analysis. To illustrate the problem of false negatives, suppose that the drug in our hypothetical actually has positive treatment effects in one of the 10 subgroups defined by age. The probability of approving the drug based on the results for the full trial population is virtually zero.⁴⁰

³⁸ Sandrine Dudoit and Mark van der Laan, *Multiple testing procedures with applications to genomics* (New York: Springer, 2008). See also the discussion in n. 22.

³⁹ We do not dispute that ex ante drug sponsors have an incentive to produce drugs that actually work. Productive drugs surely generate more revenue than unproductive ones. The problem is post hoc subgroup analysis we examine occurs after a drug sponsor has found that the average patient does not benefit from its drug. Therefore it is facing a loss equal to the cost of its drug development expenses. The only reason it has to avoid spurious correlation due to multiple testing is litigation risk from failure-to-warn products liability suits. But these suits can be foreclosed by appropriate warnings. Moreover, they do not cover the risk of a drug being less effective than alternative treatment. Finally, it is unlikely that the average ineffective drug has expected litigation costs larger than the incurred cost of development.

⁴⁰ One would need five or more subgroups that do not benefit to “show” a benefit. This is roughly equivalent to the probability that five or more successes out of ten draws from a binomial distribution with $p=0.05$. The formula and value are $\Pr(x \geq 5) = 1 - \sum_{i=0}^4 \binom{10}{i} 0.05^i 0.95^{10-i} = 0.0000275$.

Before we can suggest a compromise solution, a good question to ask is whether *post hoc* subgroup analysis can ever provide sufficiently reliable information to warrant approval of a drug, even ignoring the risks from spurious correlation and opportunistic behavior. After all, when subgroups are not specified before a trial begins, the trial is not powered – i.e., does not have sufficient sample size and thus does not generate sufficient statistical information – to estimate subgroup effects in a manner that meets the usual standards for confidence (5%) and power (20%).⁴¹ Compounding this problem is that subgroups by definition have smaller sample size than the full trial population.

We do not think, however, that these concerns render *post hoc* analysis useless. There are a number of technical reasons why such analysis may yield useful information even without a larger sample. First, sample size calculations are based on *estimates* of the variance of treatment effects in the full trial population. Because those estimates themselves have sample variation, there is a positive probability that they are in fact too high, leaving samples larger than required for accurate identification of subgroup effects. Second, because subgroups are a subset of the full trial population, they are correlated with that population. Thus analysis of the full trial population and one subgroup requires something less than a two-fold overestimate of variance to be powered to give reliable information. Third, because subgroups may be more homogenous than the full trial population, the subgroup may have smaller variation in treatment effects. This diminished variation means that doubling the number of subgroups does not double the required sample size for the trial. We shall demonstrate this in our analysis of the MGd trial.

3. Rehabilitating *post hoc* subgroup analysis

In this section we discuss two proposals that offer a compromise between (1) false positives due to opportunistic behavior and (2) false negatives or the cost of additional trials due to the FDA's cautious approach to subgroup analysis. Our aim is to extract more *reliable* information on subgroup effects from trial data that can be used to approve drugs for use in subgroups with as *little additional cost* as possible from larger sample size in the initial trial or a new trial.

Our first proposal – the use of adaptive designs – should not come as a surprise. It has been advocated by biostatisticians and regulators as a way to reduce sample size or limit harm to trial participants even in the absence of varying treatment effects across patient subgroups.⁴² The

⁴¹ See Salim Yusuf et al., Analysis and Interpretation of Treatment Effects in Subgroups of patients in Randomized Clinical Trials, 266 JAMA 93, 94 (1991).

⁴² See, e.g., Donald A. Berry, Bayesian clinical trials, 5 Nature Reviews - Drug Discovery 27 (2006) (arguing that adaptive design can be employed to lower sample size and improve treatment outcomes for enrolled patients); Scott Gottlieb, Speech before 2006 Conference on Adaptive Trial Design, Washington, DC (July 10, 2006), available at <http://www.fda.gov/oc/speeches/2006/trialdesign0710.html> (last checked Jan. 28, 2009) (observing that adaptive designs can end trials of drugs with severe side effects more quickly).

remainder of this section sketches how adaptive designs help address subgroup effects and discuss the sample size costs of those designs. The second proposal – deferred to the next section – requires a modified form of subgroup analysis to be performed by an outside consultant. It would also allow the FDA to approve a drug without the expense of further trials.

a. *What the FDA should continue to do*

Before explaining our reform proposals, we want to highlight two things that the FDA gets right in its current policy towards subgroup analysis.⁴³ First, the FDA is correct that, if the identity of sensitive subgroups is known prior to a trial, the sponsor should set the sample size for a trial so that the trial is able to estimate significant results for these subgroups.⁴⁴ For reasons mentioned earlier (correlation between subgroup outcomes and full trial population outcomes and greater homogeneity within subgroups), the additional sample size required to analyze two subgroups is not double that required to analyze the single full trial population. Therefore, the costs of powering a trial to test *a priori* specified subgroups is less than proportional to the number of groups, as is often assumed.

⁴³ However, we are concerned that the FDA makes other mistakes when aggregating information *across* clinical trials. Although the FDA may correctly apply multiple-testing corrections within trials, there is reason to be concerned that the FDA does not apply those corrections correctly across trials. Suppose a sponsor conducts an initial trial that does not show intent-to-treat effects but post hoc analysis reveals possible subgroup effects, and the sponsor conducts a second trial solely to confirm the subgroup effects. On the one hand, the second trial should be able to credit subgroup members in the first trial towards sample size requirements in the second trial. On the other hand, a significant result in the second trial may be spurious because it is itself a second test. Indeed, if one conducted 100 trials on a given subgroup, 5% would show significant effects for that subgroup even if its true effect is zero. This risk of multiple testing across trials is partly addressed by the fact that a company must inform the FDA of every trial it conducts to support an IND and by the fact that many journals will not publish an article reporting the results of a trial that has not been reported to a trial registry such as clinicaltrials.gov. See 21 C.F.R. § 312.23 (2008); Catherine De Angelis, et al., Clinical trial registration: a statement from the International Committee of Medical Journal Editors, 141 *Ann. Intern. Med.* 477 (2004). We cannot, however, find evidence from stated FDA policy or practice that the FDA understands and properly addresses these concerns. That said, it is likely the case that the extreme cost of Phase III trials limits the frequency of opportunistic behavior across multiple trials.

⁴⁴ An even better approach may be to specify subgroups not by patient characteristics at baseline, but by an algorithm that has as inputs not just those characteristics but also outcomes recorded as the trial progresses. Suppose the sponsor suspects that treatment effects may depend on one of 10 genetic markers, but is not sure which one. Instead of picking one of the those markers before the trial begins, the sponsor could, for example, specify that after ϕ fraction of subjects have enrolled, it will correlate those markers with outcomes and pick as a subgroup those subjects possessing the marker with the highest correlation with outcomes. So long as ϕ is specified before the trial begins, it is theoretically possible – though perhaps not easy – to derive a sample size to ensure this trial is properly powered. There may not be any penalty for multiple testing so that the critical p-value may remain 0.05. Nor is there a risk of opportunistic behavior by the sponsor since the FDA can implement the algorithm itself and verify the subgroup the sponsor has identified as correct.

Second, the FDA is also correct to use multiple-testing adjustments to avoid spurious results from analysis of *a priori* specified subgroups. The FDA is aware that Bonferroni adjustments may be too conservative because of the assumption that subgroups are independent. Because some patients fall in multiple subgroups or share biological features of members in other subgroups, treatment results in one subgroup may be related to effects in another subgroup. Thus, separate tests on two subgroups are less than two bites at the apple. Applying a Bonferroni adjustment in this case would result in an overcorrection for the risk of spurious correlation.

b. Adaptive trials proposal

i. Background on adaptive trials

The prototypical clinical trial is a fixed design trial. In this design, patients are randomized between a treatment group and a control group. The total number of patients enrolled in the trial – the sample size – and the fraction of patients assigned to treatment group are fixed before the trial begins and remain the same until it ends. The sample size required to run such a trial depends on the minimal size d of clinically-relevant treatment effects the sponsor wants to be able to identify and the variance σ^2 of the treatment effect from the drug.⁴⁵

The problem is that the sponsor may not know these parameters. Indeed, one of the purposes of the trial is to estimate these parameters. One solution is to use results from prior studies of the sponsor’s drug or of related drugs. When such studies are not available or are unreliable, the sponsor can use what is called an adaptive design trial. Such a trial begins without firm or completely reliable estimates of the parameters above. Instead, the trial employs real time data gathered from early-enrolling patients to refine estimates of the parameters and adjust sample size or treatment allocation based on the new estimates.

There are two types of adaptive designs that use interim data to modify sample size while the trial is in progress. In one, called a sequential-group approach, the sponsor starts with a trial that is conservatively large – using parameters at the lower end of the range for clinically-relevant effects and at the higher end of the range for variation in treatment effects – but stops the trial early if interim data suggest that treatment effects are larger than clinically relevant or have smaller variance than hypothesized. The other design, called simply an adaptive approach, does the opposite. It starts with a trial that is deliberately small and extends the trial if estimates of the treatment effect are smaller than the clinically relevant amount or estimates of the treatment effect variance are larger than hypothesized. Either adaptive design requires a larger sample size than a fixed design. Moreover, because the trial is updated after the sponsor “tests”

⁴⁵ See supra note 9.

the data by estimating treatment effects, the critical p-value may have to be reduced to account for multiple testing. The exact multiple-testing penalty has been derived in statistical literature.⁴⁶

There are also adaptive designs intended to adjust the proportion of enrolled subjects assigned to the treatment group based on interim data analysis. If, for example, outcomes in the treatment group show higher variance relative to the control group than anticipated, then the sponsor may change group assignments so that more than half of subjects get treatment. So long as the estimate of the variance of treatment effects – the difference in outcomes in the treatment and control groups – does not increase, so that the sample size remains constant, the sponsor pays no multiple-testing penalty for such an adaptive design.⁴⁷

The FDA is open to use of adaptive designs. Its Critical Path initiative begun in 2004, seeks to identify biological and statistical innovations that can improve the efficiency of clinical trials and incorporate them into the drug development and approval process. That initiative has identified adaptive designs as one area on which to focus its attention. Indeed, the FDA is expected to release a guidance on adaptive designs to clarify its thinking.⁴⁸

ii. Adaptive design for subgroup analysis

None of these adaptive designs, however, are specifically intended to address subgroup effects. They are mainly directed at optimizing over power and cost for main group effects. That does not mean that no one has thought of applying adaptive designs to estimate subgroup effects. We know of no instances, however, where the FDA has approved an adaptive design to facilitate subgroup analysis, though the FDA has considered or allowed a number of trials with adaptive designs.

How might an adaptive design be used for subgroup analysis? Consider a two-arm trial (treatment and control)⁴⁹ with sample size set to test just one hypothesis: the average treatment effect for one all enrolled patients is zero. At some interim point, the sponsor or the independent data monitoring committee (IDMC) examines the data to determine if there is a subgroup of patients on which the trial should focus because they may be particularly responsive to treatment.

⁴⁶ Cyrus R. Mehta and Nitin R. Patel, Adaptive, Group Sequential and Decision Theoretic Approaches to Sample Size Determination, 25 *Statistics in Medicine* 3250-3269 (2006).

⁴⁷ Even if the estimate of variance of treatment effects falls, the sponsor cannot stop the trial early. But if the estimate of overall variance of treatment effects rises, then the sample size increases and the sponsor must pay a multiple-testing penalty. The reason is that it was given a “real option” of testing and must pay a price for this option.

⁴⁸ Scott Gottlieb, Speech before 2006 Conference on Adaptive Trial Design, Washington, DC (July 10, 2006), available at <http://www.fda.gov/oc/speeches/2006/trialdesign0710.html> (last checked Mar. 7, 2008).

⁴⁹ This is also called a parallel-armed trial.

There are two types of data that might be used to identify subgroups: baseline characteristics measured before the start of a trial alone or treatment outcomes measured during the course of the trial. In the first case, the sponsor looks for abnormal variation in a relevant covariate. For example, if there is much more variation than expected in treatment history or in the pre-trial progression of symptoms, the full trial population can be divided into subgroups using a cut-off based on the extent of prior treatment or symptoms. In the second case, the IDMC may look at the relationship between certain covariates and treatment effects. (The IDMC is used rather than the sponsor to ensure that the sponsor does not become unblinded.) If the data suggests, for example, that certain age or ethnic subpopulations are responding better to treatment, those groups can become target subgroups for the study.

After this interim analysis, the sponsor would have to revisit the objective of the trial. There are two choices. First, the sponsor could examine just one hypothesis but limit it to a subgroup identified by interim analysis as particularly sensitive to treatment. Specifically, the null hypothesis would become: the treatment effect for *one subgroup* is zero. We assume in this case that, after the interim analysis, the sponsor would discontinue enrollment of subjects that do not belong to this subgroup, lest they waste sample size. The sponsor's other choice is to examine two or more hypotheses based on the number of subgroups discovered through interim analysis. For example, if that analysis identified two subgroups based on ethnicity, the trial might test two hypotheses: the treatment effect for whites is zero and the treatment effect for non-whites is zero.

As with adaptive designs targeting sample size adjustments, adaptive designs targeting subgroups will require a larger sample size and appropriate adjustments of the statistical methodology. These can be derived, though we do not do so here. In particular, the design may require that the sponsor pay a penalty, i.e., that the results be held to a more stringent or lower critical p-value before they are declared statistically significant, to account for the possibility of multiple testing. We explore these penalties in the Appendix.

Before we conclude our discussion of adaptive designs, it is worth noting an important weakness of these designs. Because of ethical and profit considerations, they may not be optimal for identifying side effects of treatments. If interim analysis suggests a particular subgroup may have worse side effects, both the sponsor and patient advocates will push to exclude that subgroup from further analysis. But doing so limits the amount of data we have on that subgroup and thus on the side effects of the drug.

c. Proposal for independent post hoc subgroup analysis

In this subsection we consider how it might be possible – working with a fixed, non-adaptive design trial – to use *post hoc* subgroup analysis to approve a drug without further trials. *Post hoc* subgroup analysis does not increase the risk of false positive drug approvals so long as the FDA makes appropriate multiple-testing corrections. But appropriate multiple-testing

corrections require knowledge of the number of tests the sponsor has performed. Because of the financial incentive to have its drug approved, the sponsor cannot be relied upon to truthfully report the number of tests it has performed.⁵⁰ Indeed, the FDA can be confident that the sponsor probably conducted more tests than that for which the FDA plans to adjust. Therefore, there remains a residual risk of false negatives above 5%.

i. Choosing an independent agent

A critical assumption in this logic is that the sponsor is financially interested in having the drug approved. If *post hoc* subgroup analysis were performed by a truly independent agent, then the FDA could rely upon that agent's report of the number of tests it conducted and fully eliminate the risk of false positives by means of multiple-testing adjustments. Of course the real question is whether the agent is truly independent, a topic to which we will turn in a moment.

There are two basic candidates for an independent agent: the FDA and an outside statistical consulting firm. Each has its strengths and weaknesses. The strength of using the FDA is that by doing the *post hoc* subgroup analysis itself, the FDA knows immediately the number of tests conducted. There is no need to rely on the absence of any other motive, as will be the case with an outside consulting firm. There are two weaknesses of the FDA. It has limited resources that make it difficult to maintain even the current level of scrutiny of new drug applications (NDAs).⁵¹ Moreover, the FDA is subject to political pressure. It has been criticized for being influenced both by drug companies and by political backlash following approval of unsafe drugs.⁵² These pressures are unlikely to perfectly offset to create an unbiased decisionmaker. As a result, the FDA may conduct too much subgroup analysis – at the cost of

⁵⁰ We have considered the possibility that the FDA could specify prior to a phase III trial the exact subgroups the sponsor may examine. This could be based on the subject-matter of the trial or on the FDA's knowledge of data from trials of competing drugs by other sponsors. There are two problems with this reform. First, the sponsor could specify subgroups based on subject matter as well as the FDA and has a strong financial interest in doing so. We doubt there are valuable subgroups that the FDA could propose that the sponsor will not have already considered. Second, sponsors of the competing drugs are likely to object to the FDA's use of their trial data – which is treated as a trade secret, see Article 39.3 of the Trade-Related Intellectual Property (TRIPS) agreement – in this manner. They would have a reasonable argument that this competitively favors later new drug applicants over earlier ones. It would also subtly reduce the incentive to innovate quickly.

⁵¹ See, e.g., Institute of Medicine, *The Future of Drug Safety: Promoting and Protecting the Health of the Public* 193 (2007). Drawing inspiration from PDUFA, one possible solution is to charge companies that seek drug approval for a patient subgroup to pay higher user fees to fund subgroup analysis conducted by the FDA.

⁵² See Gardiner Harris, *F.D.A. is Faulted for Drug-Safety Process*, *New York Times* (Sept. 20, 2006), available at <http://www.nytimes.com/2006/09/22/business/22fdacnd.html?ex=1316577600&en=04c9d9824b892f3b&ei=5088&partner=rssnyt&emc=rss> (last checked Jan. 30, 2009); Avery Johnson and Ron Winslow, *Drug Makers Say FDA Safety Focus Is Slowing New-Medicine Pipeline*, *Wall Street Journal* (June 30, 2008), available at http://online.wsj.com/article/SB121476772560213981.html?mod=hps_us_whats_news (last checked Jan. 30, 2009).

false positives – or too little subgroup analysis – at the cost of false negatives or more costly approval.

The alternative is an outside statistical consulting firm. Many already exist to help sponsors design and analyze data from trials.⁵³ The strength of consulting firms is, perhaps, more statistical expertise than the FDA. Unlike the FDA, which has limited resources and no need to compete, these firms have every reason to specialize and innovate because it may make it more likely they are selected to perform subgroup analysis.

The main weakness of the consulting firm approach is that these firms may not be truly independent. Sponsors are repeat players. A consulting firm may have an incentive to give a favorable analysis so as to secure repeat business from sponsors. That repeat business may be for subgroup analysis or some other statistical service. This is a lesson well learned from the corporate accounting scandals from earlier this decade.⁵⁴ Perhaps the indirect influence of sponsors can be addressed by requiring the FDA to select the outside consultant to perform *post hoc* analysis, by blinding the sponsor to the outside firm selected, and by banning firms that perform *post hoc* analysis from providing other statistical services to sponsors. We wonder, however, whether the agency will always be able to keep the identity of the consulting firm secret, even after the analysis is completed and the FDA has made its regulatory decision concerning a sponsor's drug. Moreover, restricting the consulting firms' scope of business will limit their ability to attract talent and incentive to innovate in the area of subgroup analysis since it comes at the cost of other lines of business.

A second weakness of using consulting firms is that “independence of the sponsor” is not the same thing as “motivated to reduce false positives.” True independence only guarantees the consulting firm will not be swayed by the profit interests of the sponsor. It does not guarantee that the consulting firm extracts the most and reliable data from *post hoc* analysis after it is chosen to perform that analysis. This problem is one which economists call moral hazard. Independence merely substitutes the sponsor's interests with those of the consulting firm. Most likely this is cost minimization, which may imply too many false negatives or false positives, whichever minimize the consulting firm's labor expense.⁵⁵

ii. Statistical methods to guarantee independence

⁵³ See, e.g., Cytel Statistical Software and Services, founded by Cyrus R. Mehta and Nitin R. Patel, and Target Analytics, Inc., run by Mark van der Laan.

⁵⁴ See Joel S. Demski, Corporate Conflicts of Interest, 17 J. Econ. Persp. 51, 57 (2003).

⁵⁵ The outside consulting firm must also be concerned about not doing too many tests. Each useless test it performs increases the multiple-testing adjustment for any positive finding. Minimizing false negatives requires internalizing this negative externality. Since false negatives are unobservable, the FDA cannot directly incentivize the consulting firm to do so. And the FDA certainly should not give the firm an incentive keyed to drug approval, because then it would have incentives like the sponsor and replace false negatives with false positives.

To address the problem that neither the FDA nor the outside consultant may be truly independent of the sponsor, we propose two statistical methods to limit either agent's ability to skew the analysis in favor of the sponsor.⁵⁶ For convenience, we shall speak as if the consulting firm has been chosen to conduct the analysis.

No-outcome data analysis. The first approach would provide the consultant with all the data from the trial *except variables that identify treatment assignments and health outcomes* and ask it to identify subgroups based on baseline characteristics that exhibit “remarkable and relevant variation” in the trial data. (This is similar to one of the approaches used to identify subgroups for the adaptive design trials discussed in the last subsection.) The consultant would not be asked to perform the *post hoc* subgroup analysis; that could be conducted by the sponsor, though the FDA would rely upon positive treatment effects only for the subgroups identified by the consultant. Whatever positive subgroup results the sponsor reports, the FDA would apply a multiple-testing adjustment based on all the subgroups reported by the outside consultant.

In order to identify remarkable variation, the consultant needs to have a sense of what normal variation would be. It could estimate normal variation in baseline characteristics from trials of the sponsor's drug or prior studies in the literature. The consultant would have to be sensitive to exclusion and inclusion criteria, which can affect the applicability of prior data to the current trial sample. Moreover, the consultant would have to keep in mind that any subgroup it identifies should be defined by variables that are plausibly relevant (from our current biological understanding of the disease targeted by the sponsor's drug and the pharmacology of that drug) to the treatment effects of the drug.

Split-sample analysis. The second statistical method we propose to ensure that the consultant's choice of subgroups is not influenced by the drug company requires splitting the data from a trial into two parts. One part would be called the exploratory subsample and the other part the confirmatory subsample. Importantly, the FDA must split the sample to ensure the drug company has no influence and the sample should be split randomly to ensure the samples are statistically independent. The consultant would only be given the exploratory subsample and be asked to conduct a full *post hoc* subgroup analysis on that subsample to identify subgroups that

⁵⁶ These methods do not address other problems, such as the limited resources of the FDA or the insufficient motivation of outside consulting firms. If the statistical methods we discuss help ensure that the consultant truly cannot manipulate the data to increase false positives, then one might address the problem of a consultant's motivation by giving it stock in the sponsor. (This is identical to extracting the outcome or a random subsample of data from the data archives of the drug sponsor. We consider granting the consultant stock instead because it is virtually impossible to separate the sponsor from knowledge of its data.) We do not advocate this because it is too radical and would be politically infeasible. That said, giving the consultant some sponsor stock is not the same as allowing the sponsor to conduct the entire *post hoc* subgroup analysis because the statistical methods we propose in the main text require that the consultant not have access to certain data that the sponsor already has, or could easily obtain.

respond better to the drug.⁵⁷ The sponsor would then be allowed to perform *post hoc* subgroup analysis on the confirmatory sample using only the subgroups identified from the exploratory subsample by the consultant. As before, the FDA would apply a multiple-testing penalty based on all the subgroups reported by the outside consultant. If, after such penalty, the confirmatory subsample validates the positive subgroup effects from the exploratory subsample, the FDA could approve the drug only for those subgroups.

Both statistical methods ensure that subgroups are identified independent of the interests of the sponsor. Since the first method does not give the consultant access to outcome data, it cannot choose subgroups to help or hinder the sponsor. Since the second method requires the sponsor to limit its subgroup analysis to a subsample that is statistically independent of the subsample analyzed by the consultant, the consultant's analysis cannot help the sponsor engage in data mining. Moreover, neither method requires that the FDA impose any additional multiple-testing penalty beyond one based on the total number of subgroups identified by the consultant.

Each statistical method also has its shortcomings. The weakness of the no-outcome data approach is that the subgroups with the most remarkable variation may not be perfectly correlated with the subgroups that have positive and significant treatment effects. Abnormal variation is just one factor that suggests differential treatment effects; it does not guarantee them. The main concern with split-sample approach is that the *post hoc* subgroup analysis, which would be underpowered even if performed on the whole trial sample, is particularly underpowered if performed on subsamples. This will increase the risk of false negatives. This risk may be considered the cost of independence under this method. In short, the two statistical algorithms reduce, but do not eliminate false negatives.

4. An illustration with motexafin gadolinium

In this section we illustrate our two statistical algorithms for ensuring independent *post hoc* subgroup analysis by applying them to a real world example: motexafin gadolinium (MGd) for patients with brain metastases from solid lung tumor. MGd is sponsored by Pharmacyclics (ticker PYCY), a small biotech company that branded the drug as Xcytrin. We first provide some background on clinical testing of the drug and then discuss *post hoc* subgroup analysis of the testing results.

⁵⁷ The sponsor could not be asked to do this because it would likely be able to derive the confirmatory subsample from the exploratory subsample and the full sample, which it already possesses. This would allow it to choose subgroups ostensibly on the exploratory subsample but truly on the full sample. The result would be almost the same as *post hoc* subgroup analysis by the sponsor.

a. Background on MGd

Tumorous cancers in one part of the body often spread – or metastasize – to other parts of the body. In up to 24% of all cancer patients, they spread to the brain.⁵⁸ The risk is especially severe with lung cancer, where up to 50% of patients experience brain metastasis⁵⁹ and metastasis occurs earlier than with other cancers.⁶⁰ Most patients with brain metastases die. Median survival on whole brain radiation therapy, the typical conventional treatment, is only 4 months. For those who do manage to survive, however, there is a major risk of neurological impairment.⁶¹

MGd is a drug that demonstrated the ability to increase the radiation response of tumor cells in preclinical studies. Pharmacyclics sought to market the drug as a treatment for brain metastases. The company filed an investigational new drug (IND) application with the FDA to begin clinical testing of the drug in human patients. After a successful Phase I/II study,⁶² the company began a Phase III study (called trial 9801) that enrolled patients with any type of cancerous tumor who developed brain metastases. Subjects were randomized to either whole brain radiation therapy (WBRT) alone (the control arm) or MGd and WBRT (treatment arm). Unfortunately, this study did not find a statistically significant treatment effect with respect to median survival or time to neurological impairment.⁶³

⁵⁸ J.B. Posner, *Neurological complications of cancer* (Philadelphia: F.A. Davis Comp., 1995).

⁵⁹ M. Stuschke, W. Eberhardt, C. Pottgen, et al., Prophylactic cranial irradiation in locally advanced non-small-cell lung cancer after multimodality treatment: Long-term follow-up and investigations of late neuropsychological effects, 17 *J. Clin. Oncol.* 2700-2709 (1999); T.J. Robnett, M. Machtay, J.P. Stevenson, et al., Factor affecting the risk of brain metastases after definitive chemoradiation for locally advanced non-small-cell lung carcinoma, 19 *J. Clin. Oncol.* 1344-1349 (2001).

⁶⁰ Posner, *supra* note 58.

⁶¹ Minesh P. Mehta, Patrick Rodrigus, C.H.J. Terhaard, et al., Survival and Neurological Outcomes in a Randomized Trial of MotexafinGadolinium and Whole-Brain Radiation Therapy in Brain Metastases, 21 *J. Clin. Oncol.* 2529 (2003).

⁶² The sing-armed study found a 72% radiologic response rate. P. Carde, R. Timmerman, M.P. Mehta, et al., Multicenter phase Ib/II trial of the radiation enhancer motexafin gadolinium in patients with brain metastases, 19 *J. Clin. Oncol.* 2074-2083 (2001).

⁶³ Median survival was 5.2 mo. on treatment versus 4.9 mo. on control (p = 0.48). Median time to impairment of neurological function was 9.5 mo. on treatment versus 8.3 mo. on control (p = 0.95). Mehta et al., *supra* note 61, at 2533 (Fig. 2, panel C).

One bright spot in trial 9801, however, was that patients with specifically lung cancer did experience statistically significant extension of time to neurological impairment.⁶⁴ So Pharmacyclics conducted a second Phase III trial (called Trial 0211) targeting only lung cancer patients. Unfortunately, this second trial was unable to validate the results from the initial trial. This is illustrated in Table 2, which summarizes the results from trial 0211 and from lung cancer patients in trial 9801. According to the first panel, whereas the relative hazard rate for neurological impairment⁶⁵ was 0.61 ($p = 0.05$) in the initial trial, it was merely 0.78 and not significantly different from 1 ($p = 0.1$) in the second trial.

Trying to explain the discrepancy between the trials and to salvage MgD for a new drug application (NDA), Pharmacyclics conducted a *post hoc* subgroup analysis. According to the company, this analysis revealed a problem at some of the study centers in France. Although the trial protocol required that subjects be randomized to treatment as soon as they were diagnosed with brain metastases, the French centers waited several weeks or more after diagnosis before randomizing subjects to treatment.⁶⁶ In the interim, the centers gave subjects chemotherapy⁶⁷ (hence these subjects are labeled “controlled” patients in the data). Moreover, subjects were ultimately randomized in trial 0211 only if their brain tumors progressed despite chemotherapy, i.e., their tumors were resistant to treatment.⁶⁸ So this was a self-selected group of tumors.

⁶⁴ The median patient on WBRT and MgD did not experience neurological impairment in 24 months while the median patient on WBRT alone experienced impairment at 7.4 mo ($p = 0.048$). Mehta et al., *supra* note 61, at 2533 (Fig. 2, panel C).

⁶⁵ Neurological impairment was judged by a battery of standardized neurocognitive tests. The tests were scored by blinded graders. Patients were said to be impaired if the composite score was at least 1.5 standard deviations worse than the mean of the test’s age-adjusted distribution. Christina A. Meyers, et al., *Neurocognitive Function and Progression in Patients With Brain Metastases Treated With Whole-Brain Radiation and Motexafin Gadolinium: Results of a Randomized Phase III Trial*, 22 *J. Clin. Oncology* 157, 158 (2004). The hazard rate for impairment is the rate at which patients are judged impaired, i.e., it is the fraction of additional patients judged impaired each month. The relative hazard rate is the ratio of the hazard rate in the treatment group to the rate in the control group.

⁶⁶ Ordinarily, a drug sponsor is responsible for ensuring that its study centers follow trial protocols. The company attributes the problem in this case to recruitment difficulties. Whereas the trial 9801 had 401 total patients, of which only 251 had lung cancer, trial 0211 required 554 lung cancer patients to have the statistical power to validate the positive results for lung cancer patients from the initial trial. This required the second trial to recruit patients from 90 treatment centers throughout the world, more than double the 40 centers involved in the initial Phase III trial. The company argues that it is difficult to precisely enforce the protocol with so many centers involved in a study. Personal communication with Richard Miller, former CEO of Pharmacyclics, Mar. 14, 2008.

⁶⁷ Chemotherapy is not thought to be a reliable treatment for brain metastases because chemotherapy relies upon drug delivered by blood and the brain tumor is somewhat protected from chemotherapy drugs by the blood-brain barrier.

⁶⁸ The implicit but reasonable assumption here is that a brain tumor resistant to chemotherapy is resistant to any other form of treatment.

Whereas “uncontrolled” subjects randomized immediately had a mix of resistant and non-resistant brain tumors, the controlled subjects ultimately randomized in the problematic French centers largely had tumors resistant to WBRT. This placed MGd at a disadvantage in these centers.⁶⁹

Not surprisingly, these subjects also did not show benefit from MGd. Excluding these late-randomizing centers from the analysis revealed that drug had a statistically significant effect on delay until onset of neurological impairment. As reported in the second of panel of Table 2, the relative hazard rate for uncontrolled subjects on MGd in the 0211 trial was 0.56 ($p = 0.02$). In other words, while MGd proved effective among subjects with uncontrolled brain tumor, this effect was masked by including subjects with controlled brain tumor in the main analysis. Pharmacyclics filed an NDA with the FDA relying on this subgroup analysis. But the FDA did not credit the company’s explanation and finally rejected its NDA in December 2007.⁷⁰

Of course, this is the company’s explanation for its unsuccessful final Phase III trial and it had a financial stake in getting MGd approved for some subgroup. Our aim is to scrutinize these claims by taking the role of an outside consultant and checking whether analysis of data in a manner that is independent of the financial interests of the sponsor identifies the same sensitive subgroups that the sponsor identified, namely the subjects with uncontrolled tumor.

b. No-outcome data analysis

Our first analysis examines the final Phase III study (trial 0211) data stripped of outcome variables. The idea is that an outside consultant without outcome data would not be able to select subgroups that would financially benefit the company because it does not know whether any subgroups had better or worse outcomes than average among the full trial population. Instead, this consultant would identify subgroups by searching for baseline characteristics on which current trial subjects had excess variation relative subjects in previous trials or in the population. These characteristics could then be used to define subgroups on which the drug company could perform *post hoc* subgroup analysis. If and only if that subgroup analysis suggested the subgroup responded positively to the drug, the FDA should approve the drug for use in that subgroup.

⁶⁹ The selection story is a bit more complicated. The delay also screened out patients who had died from, inter alia, the brain metastases prior to randomization. Since early death is an indicator of a more severe brain tumor, this mortality screen likely selected for less severe brain tumors. It is probably the case that the selection on the basis of resistance to chemotherapy (which likely reduced the effect of MGd) was more significant than selection based on survival (which possibly increased the effect of MGd). The reason is that median survival following diagnosis with brain metastases is 4 months, so it is unlikely that mortality was a material screen in the first two weeks following diagnosis. Yet it is in these first two weeks that the company found a significant delay in neurological impairment among patients treated with MGd.

⁷⁰ See supra note 23.

To implement this algorithm, we compare the variation of certain medical characteristics in the trial 0211 sample with variation of those variables in the initial Phase III study (trial 9801) sample. The variables include features of the primary (lung) tumor, treatment of the primary tumor, features of the brain metastases, the treatment of the brain metastases before enrollment, and neurological impairment at baseline. Table 3 reports the ratio of variances of each variable across the two samples and the p-value for the hypothesis test that the variances are equal across the samples after adjusting for multiple testing.

Eight subgroups stand out. The 0211 trial had excess variation in the variables: days from diagnosis of brain metastases to randomization, extracranial metastases, baseline Trail B score, and whether the study center was in Canada. The delay variable captures some of the company's concern that patients in some French centers received chemotherapy for a few weeks before being randomized and that the brain masses that survived this chemotherapy were more resilient. The trail B is a test of cognitive ability where the subject is asked to follow a "trail" on a sheet of paper with his pencil. A lower score is better: it indicates less time was required to follow the trail. It also indicates that the brain metastasis probably has not advanced beyond the point where it can be treated.

The 0211 trial also had insufficient variation in: primary tumor (PT) resected without recurrence, primary tumor is large cell carcinoma, primary tumor controlled, and study center is in the US. The PT resected variable is one of five categories into which a primary tumor is categorized at the time a patient is randomized.⁷¹ The resected category indicates that the primary tumor was surgically and successfully treated and the patients only remaining concern is the brain metastases. The PT controlled variable is the complement of the variable that the company identified as being responsible for the failure of the 0211 trial.

Having identified eight subgroups with remarkable variation, we now check to see if MGd was particularly effective amongst these subgroups in the full 0211 data. For binary variables, subgroups are defined by their two states. For continuous variables subgroups in the 0211 data are defined by whether a characteristic lies above or below the median for that characteristic in the 9801 data. Table 4 summarizes our subgroup analyses with a multiple-testing adjustment that accounts for 16 tests (eight variables with excess variation and two subgroups for each variable). We find four subgroups have significant treatment effects, i.e., lower rate of neurological impairment: subjects with little delay before randomization, subjects at U.S. study centers, subjects with low (good) baseline Trail B scores, and subjects with uncontrolled brain metastases.

⁷¹ The tumor may have been (1) "newly diagnosed," which means the primary tumor and the brain metastasis was diagnosed at the same time, (2) surgically removed or resected without recurrence, (3) treated for less than 4 weeks without clear indication that it has progressed, (4) treated for greater than 4 weeks with no sign of progression, or (5) treated for any amount of time with evidence of progression.

Thus the no-outcome data analysis would support Pharmacyclics' case for approval for subjects with uncontrolled brain tumors. The other subgroups that survive no-outcome data analysis are consistent with the company's theory that MGd works on less resistant tumors. Subjects randomized quickly and subjects in the U.S. are less likely to have received chemotherapy before randomization. As for subjects with low trail B scores, these are subjects whose tumors are not so far along that they already seriously impede subjects' cognitive capacity. It makes sense the treatment is also likely to work in these cases.

c. Split-sample analysis

Our second (and preferred) proposal is to split the data from trial 0211 into two subsamples, have an outside consultant to identify subgroups via *post hoc* subgroup analysis on one (exploratory) sample, and then allow the drug sponsor to validate significant treatment effects for the other (confirmatory) sample. Only if a subgroup identified by the outside consultant demonstrates statistically significant effects – after a multiple-testing penalty – in the confirmatory sample should the FDA approve the drug for that subgroup. Because the outside consultant does not have access to the confirmatory sample and because the FDA would only credit evidence of treatment effects from the confirmatory subsample, the consultant cannot rig its analysis to help the drug sponsor.

We begin our simulation of the split-sample analysis by randomly dividing the 0211 trial sample into a 20% exploratory sample and an 80% confirmatory sample.⁷² The second step is to conduct a subgroup analysis of the exploratory sample where subgroups are defined according to baseline variables. The results of this analysis are presented in Table 5. The vertical panel labeled “Value A” gives the relative hazard of neurological impairment and the p-value for each variable at value 0 for binary variables and at the first quartile for continuous variables. The vertical panel labeled “Value B” gives the relative hazard of neurological impairment and the p-value for each variable at value 1 for binary variables and at the third quartile for continuous variables.⁷³ None of the p-values have been adjusted for multiple testing.⁷⁴ We choose a subgroup (now defined both by a given variable and a specific value for that variable) for the third step in our simulation if it has a p-value less than 0.05, i.e., if we can be 95% confident that

⁷² Our selection of a 20-80 split is arbitrary. Further statistical analysis is required to determine the optimal split of the sample. The larger is the exploratory sample, the greater is the probability of identifying a subgroup that benefits but the smaller is the probability that one will be able to confirm that it benefits.

⁷³ The first two columns of data present the ratio of relative hazards among the two subgroups defined for each variable and the p-value for this ratio.

⁷⁴ No multiple testing penalty is required when analyzing the exploratory sample. The only purpose in that sample is to identify certain subgroups, relative to other subgroups, that have a better response. Moreover, the size of the exploratory sample is too small to pass any of the usual statistical tests (e.g., $p = 0.05$), let alone ones that adjust for spurious correlation from multiple testing.

the membership in the subgroup improves treatment effects. The two subgroups we identify in this manner are (1) enrollment at center other than one in France and (2) not having previously received the chemotherapy drug carboplatin. The relative hazard rate (into neurological impairment) for subjects outside France and on MGd is 0.32 ($p = 0.05363$) and for subjects not having received carboplatin and on MGD is 0.039 ($p = 0.07833$).

The last step in the split sample analysis is to estimate the treatment effects for these subgroups in the confirmatory sample. This analysis reveals that the relative hazard rate for MGd among subjects outside France is 0.676 (raw $p = 0.038$, adj. $p = 0.076$) and among subjects who had not previously been treated with carboplatin is 0.905 (raw $p = 0.58$, adj. $p = 0.58$). Even with the high correlation ($\rho = 0.39$) in membership across these subgroups, the effect of MGd outside France is not statistically significant after adjusting p-values for multiple tests (on two groups). The effect of type of prior chemotherapy is not significant even without the multiple-testing adjustment.

So it appears that our spit sample analysis – unlike the no-outcome data analysis – fails to support Pharmacyclics explanation for why the 0211 trial did not validate the 9801 trial. To be precise, however, all we have shown is that the particular sample split we randomly drew chose did not validate Pharmacyclics claim. Perhaps another random split would validate their claim. Indeed, a better way to characterize the value of the split-sample analysis for eliminating alleged false negatives is to ask, what fraction of splits would validate Pharmacyclics’ claim that MGd works in the subgroup of patients whose brain tumor was not controlled via chemotherapy?

To conduct this analysis we drew 100 splits of the 0211 trial data and ranked the subgroups in the exploratory stage in order of statistical significance of relative hazard for neurological impairment, just as Table 5 did for our initial split. The first row of Table 6 provides the distribution of p-values that the “uncontrolled” subgroup takes across the hundred splits, with a smaller p-value indicating a larger significance of that difference. We find that in 35% of draws uncontrolled has a p-value of less than 0.05. Thus in 35% of cases, the “uncontrolled” subgroup advances to the validation stage. Further, in 31% (= 11/35) of these cases, the effect among the uncontrolled subgroup is validated in the confirmatory sample after a conservative Bonferroni multiple-testing adjustment that accounts for the number of subgroups that emerge from the exploratory analysis in each sample split. See row 3 of Table 6. In other words, in only 11 % of cases, the split-sample analysis confirms Pharmacyclics’ claim that MGd works so long as the patient’s brain tumor is not previously treated with chemotherapy.

Because in the case of MGd we have not one but two Phase III trials, there is one other test we can do – in the spirit of split sample analysis – to verify Pharmacyclics’ claim about uncontrolled patients. We can check whether the subgroup effects identified by the company, by the no-outcome data analysis, and by 11 % of the split-sample analyses above can be validated

by in the subsample of lung cancer patients in the 9801 trial.⁷⁵ As can be seen in Table 2, the uncontrolled subgroup is indeed associated with statistically significant treatment effect (RH = 0.48, $p = 0.03$).

Conclusion

Roughly one in five drugs that enter clinical testing fails to prove that it is effective and safe. Even in phase III, the failure rate is 36%.⁷⁶ The cost of failure is a higher cost of developing drugs. By one commonly cited estimate, whereas the expected cost of clinical testing per drug that begins human trials is \$316.3 million, the cost of clinical testing after accounting for the risk of failure – i.e., the cost per drug that is approved – is \$802 million.⁷⁷

Sometimes failure is just that: the drug has no value. But other times a drug is right for some patients and wrong for others. Denying approval for a drug that benefits some patients, but not the average patient, increases the costs of drug development but not the benefits. Ideally, one would like to salvage such a drug by allowing its use for the non-average patient who would benefit.

The challenge is bad behavior or moral hazard by drug sponsors. With *post hoc* subgroup analysis (or data dredging in less polite language), sponsors will nearly always be able to find a subgroup of patients who appear to benefit from a drug. Currently, the FDA addresses this risk of spurious correlation by requiring sponsors to validate their findings with additional clinical trials. But this may cost tens of millions of dollars, and in turn increase the price of drugs.

This paper offers a combination of institutional designs and statistical methods that can limit the risk of spurious findings – or false positives – from *post hoc* subgroup analysis without requiring additional, whole trials. Our proposal for adaptive trials allows the use of subgroups to revise the hypothesis tests in a trial with little additional sample size. Our proposal for independent statistical analysis, when combined with subgroup analysis without outcome data or subgroup analysis validated on a split sample, can actually identify subgroups increasing the risk of false positives or requiring additional sample size. In other words, our proposals offer an

⁷⁵ Indeed, we could ask the same thing of any subgroup identified by the no-outcome data analysis or the split-sample analysis that was verified in the 0211 analysis. Though we would then have to apply a multiple-testing adjustment accounting for all the groups we test. Since only uncontrolled survives our 0211 validation, that is all we test.

⁷⁶ Roughly 90 percent of drugs submitted in NDAs are approved. Christopher P. Adams and Van V. Brantner, Estimating The Cost Of New Drug Development: Is It Really \$802 Million?, 25 Health Aff. 420, 422 (Exhibit 1) (2006).

⁷⁷ Joseph A. DiMasi, Ronald W. Hansen, Henry G. Grabowski, 22 J. Health Econ. 151, 167 (2003).

approach to reduce the rate of failure in clinical trials, without a higher risk of false positives or with minimal additional clinical testing costs.

While our proposals may be helpful, we recognize they are not panaceas. It is possible that a drug which fails the average patient standard used by the FDA may not in fact be helpful to any subgroup of patients. Our methods for identifying false negatives – drugs that have value for a subgroup of patients but not the average patient – may not identify every false negative. Finally, even if a drug is approved for the right subgroup, doctors may use it for the wrong subgroups or use it off label. These are forms of false positives that we cannot address. Indeed, by increasing the number of drugs available to doctors, we increase the risk of post-approval false positives.

The reason we believe our proposals are worth pursuing, however, is that the alternatives – using an average-patient standard or requiring additional trials – are worse. As we explained in Section 1, the average-patient standard implicitly assumes that doctor *always* give the drug to the wrong subgroup. While doctors may not be perfect at sorting patients to drugs, we do not believe they are as bad as the FDA’s standard assumes. Moreover, trials – especially Phase III trials – are very expensive.⁷⁸ Trials focusing on subgroups are even more costly. Because fewer patients are members of the subgroup than the full trial population, a trial focusing on a subgroup will take longer complete recruitment. This, in turn, increases the opportunity costs of the trial.

It is natural to wonder whether our proposal to eliminate false negatives in drug approval can also help eliminate false positives. That is, should our proposals be used to identify drugs that have a no effect or a side effect for a patient subgroup, even though they are effective and safe for the average patient and thus FDA-approvable? The answer is complicated. As we have mentioned, adaptive trials are not helpful for identifying side effects.⁷⁹ With respect to our other proposals, the answer depends on whether there is reason to believe that the FDA has a skewed incentive to disapprove helpful drugs. The central problem that motivates our proposals is moral hazard by drug sponsors. Sponsors have a financial incentive to find subgroups that show benefits from a drug whether or not the subgroups actually benefit from the drug. Unless there is a corresponding incentive on the part of the FDA to disapprove a drug for certain patients though it does not harm them, there is no reason why the FDA cannot itself conduct *post hoc* subgroup analysis to identify subgroups that do not benefit from a drug. Since the FDA is conducting the analysis, it knows the number of subgroups it has tested and thus can apply a multiple-testing

⁷⁸ All phase III testing – usually two trials – is estimated to cost \$205 million per drug than enters phase III. See DiMasi, *supra* note 26, at 162.

⁷⁹ Adaptive trials use data from early enrollees to identify subgroups and then modify the trial to enroll more patients from those subgroups so as to explore drug effects in those subgroups. When identifying subgroups with side effects, this will require enrolling more people who one suspects will experience side effects. This is unethical and unlikely to be approved by an Institutional Review Board. See the last paragraph of Section 3.b.ii.

penalty to its own analysis. The only reform that we can unambiguously recommend to address false negatives is that the FDA be sure to apply multiple testing penalties when it conducts or requires that drug sponsors check side effects in certain subgroups.⁸⁰ Just as data dredging allows a drug company nearly always to find some subgroup of patients who appear to benefit from a drug, the more subgroups the FDA tests for side effects, the more likely the agency will find side effects when they do not actually exist.

Finally, we are aware that our paper raises a number of statistical questions, the answers to which would help the FDA refine regulations that allow certain *post hoc* subgroup analyses to inform approval decisions. For example: What is the proper multiple testing penalty for adaptive design trials? Can sponsors use moments other than the variance, such as the mean or skew, to identify subgroups in the no-outcome data analyses? What are the appropriate proportions (20-80 or something different) to use when dividing a sample into an exploratory and a confirmatory subsample? Should the multiple-testing penalty applied to tests on the confirmatory subsample account for the fact that that subsample is smaller than the full sample? Under what conditions will the no-outcome data analysis eliminate more false negatives than the split-sample analysis. We leave these questions for future analysis by statisticians.

⁸⁰ A subtle point related to this reform is that the FDA ought not blindly to apply multiple testing penalties to subgroups that a drug sponsor voluntarily offers up as experiencing side effects. The reason is that the sponsor has an incentive to report that it tested far more subgroups than it actually did so as to raise the multiple testing penalty for the subgroup that it does identify as having a side effect. In a backhanded way, this allow the drug sponsor to look responsible – it offered up a subgroup that should not get its drug – but not actually lose any sales – the multiple testing penalty would render the result for that subgroup insignificant. A better solution is for the FDA by itself to conduct *post hoc* subgroup analysis designed to identify subgroups who do not benefit from a drug.

Appendix

The table below summarizes the four basic options in a subgroup-identifying adaptive design and our speculation as to the appropriate multiple-testing penalty. The rows indicate whether interim analysis employed outcome data or not. The columns indicate whether the sponsor added hypothesis tests after the interim analysis.

Table 1. Multiple testing penalties in adaptive trials.

Data used to identify subgroups	Number of additional hypothesis tests added to study	
	Zero	One or more
Covariates (not outcomes)	No penalty	Penalty for adding one or more hypothesis
Outcomes	Must keep other subgroups in evaluated population, plus a penalty for using outcome data	Penalty for adding second hypothesis, must keep other subgroups in evaluated population, plus pay a penalty for using outcome data

If outcome data are not used to identify subgroups and no additional hypothesis tests are added to the study, then there is no need to impose a multiple testing penalty, so long as the trial must proceed until the sample size specified prior to starting the trial is achieved. The reason is there was no testing of treatment effects in the interim analysis and the number of tests remain the same as when the trial began. Even though the sponsor may choose a subgroup with low variance with respect to covariate characteristics, so long as the data employed in the interim analysis (patients' baseline characteristics) are unrelated to the data relevant for estimation of the treatment effect (patients' treatment assignment and health outcomes), there is in essence additional testing of treatment effects in the interim analysis. The general idea is, so long as one analyzes a subset of the final data set that contains no information about treatment effects and does not increase the *number* hypotheses to be tested, there is no multiple testing penalty for changing the nature of the hypothesis to be tested with the final data.

If outcome data were used to identify subgroups, there should be a multiple testing penalty even if no additional hypothesis tests were added. The reason is that the sponsor was able to test whether treatment effects are significantly positive for a subgroup during the interim analysis. Even with that subset of the final sample, it is highly likely that data mining would uncover at least one subgroup with significant treatment effects in the subsample. As a result, the sponsor would have been given the option to change the hypothesis test's scope based on treatment effects. It must pay a price for that.

This price is difficult to calculate since we may not know how many tests were performed to identify a subgroup. It helps if that the IDMC conducts the interim analysis because it has less incentive to engage in data mining and in any case is more likely truthfully to report the number of tests performed. But if the sponsor has a role on that committee or the

IDMC is not otherwise truly independent, institutional design may not help with calculating the multiple-testing adjustment.⁸¹

In the cases where the sponsor adds one or more hypothesis to the study, it must pay an additional price for multiple testing on top of the price it pays based on the data employed to conduct the interim analysis. The reason for this penalty is obvious – the number of hypothesis test has increased – and the size of the incremental penalty is straightforward to calculate.

⁸¹ In that case, we speculate – though have not confirmed – that requiring the sponsor to include the excluded groups along with the newly targeted subgroup in the final analysis may address the problem that the FDA may not know the number of tests performed in the interim analysis. Suppose, for example, that interim analysis after 10% of the sample is enrolled reveals that only young patients have a significant treatment effect. We recommend, when the sponsor tests the hypothesis that the treatment effect among young patients in its final empirical analysis, that the sponsor be required to use the entire sample and not just young subjects. Specifically, the sample tested should include the elderly patients from the initial 10% sample even though they are not nominally the subject of the hypothesis test. Our crude logic is that the larger the number of tests performed, the worse the relative performance of subgroups excluded from the modified hypothesis test, and the larger is the cost or penalty to the sponsor of having to include the excluded subgroups in the final empirical analysis. Including the elderly from the 10% sample automatically make it less likely that the sponsor will be able to show that the drug works among young patients.

Tables and figures

Table 2. Subgroup treatment effects in trial 0211 and trial 9801.

Group	Trial 0211			Trial 9801		
	n	RH	raw p-value	n	RH	raw p-value
All	554	0.78	0.1	251	0.61	0.05
PT controlled?						
Yes	140	1.71	0.059	93	0.94	0.86
No	414	0.56	0.002	158	0.48	0.03
Newly diagnosed?						
Yes	259	0.59	0.032	109	0.47	0.046
No	295	0.92	0.69	142	0.74	0.37
Time from BM to Tx						
Tx ≤ 2 wks	274	0.6	0.022	119	0.78	0.5
2 < Tx ≤ 4 wks	161	0.78	0.41	69	0.63	0.29
Tx > 4 wks	119	1.23	0.5	63	0.33	0.09
Prior chemotherapy						
No	315	0.67	0.06	155	0.57	0.07
Yes	239	0.91	0.66	96	0.72	0.42
Trail B score						
Low	258	0.53	0.012	121	0.76	0.44
High	254	1.02	0.92	96	0.7	0.41
Country						
USA	185	0.39	0.0048	123	0.76	0.45
Netherlands	11	5.6	0.14	54	0.38	0.074
Canada	163	0.72	0.26	46	0.4	0.14
UK	0			21	1.11	0.89
France	117	1.49	0.21	7	0.82	0.89
Germany	47	0.61	0.26	0		

Notes. RH = relative hazard for MGd plus whole brain radiation therapy (WBRT) versus WBRT. Raw p-value does not adjust for multiple testing. Trial 9801 was original Phase III trial. Trial 0211 was second Phase III trial. PT = primary tumor. BM = brain metastases. Tx = treatment in treatment or control group.

Table 3. Identification of subgroups in non-outcome analysis: ratio of variance for baseline characteristics in 0211 trial versus 9801 trial.

Covariate	Ratio	p-value
Extracranial metastases?	1.73	0.0006
USA?	0.89	0.0006
PT controlled?	0.81	0.0025
Canada?	1.39	0.0068
PT resected without recurrence?	0.50	0.0137
Baseline Trail B score	1.63	0.0158
PT is large cell carcinoma?	0.60	0.0267
Days from diagnosis to randomization	5.29	0.0267
PT treated, <=1 month follow-up?	1.72	0.0554
PT is non-small-cell carcinoma?	0.80	0.1200
RPA status	0.71	0.1519
Baseline weight	1.27	0.2171
Sex	0.98	0.2171
PT has squamous histology?	0.78	0.2249
PT treated, >1 month follow-up?	0.85	0.2249
PT treated and progressing?	1.24	0.2734
PT has other histology?	1.90	0.3143
Prior chemotherapy?	1.04	0.4045
Baseline delay score	0.90	0.4045
Baseline COWA score	0.88	0.4045
Baseline Trail A score	0.59	0.4045
Baseline height	1.11	0.4075
Multiple BM lesions?	0.89	0.4154
PT newly diagnosed and/or untreated?	1.01	0.5843
Age 65+?	0.96	0.6924
Karnofsky Performance Score >= 90?	1.01	0.7010
PT is adenocarcinoma?	0.99	0.8469
Caucasian?	1.07	0.8469
Baseline recall score	0.97	0.8469
Baseline recognition score	0.98	0.9619
Other race?	1.00	0.9931

Notes. P-value adjusts for multiple testing. PT = primary tumor (i.e., lung cancer). BM = brain metastasis.

Table 4. Treatment effect in 0211 trial for subgroups identified by no-outcome data analysis.

	Relative hazard	p-value
Canada? No	0.82	0.3782
Canada? Yes	0.72	0.3782
Extracranial metastases? - Low	0.85	0.5467
Extracranial metastases? - High	0.71	0.2800
PT is large cell carcinoma? No	0.78	0.2800
PT is large cell carcinoma? Yes	0.90	0.8853
Days from diagnosis to randomization - Low	0.55	0.0325
Days from diagnosis to randomization - High	1.14	0.6646
USA? No	0.99	0.9500
USA? Yes	0.39	0.0325
Baseline Trail B score - Low	0.50	0.0480
Baseline Trail B score - High	0.96	0.8853
PT controlled? No	0.56	0.0320
PT controlled? Yes	1.71	0.1573
PT resected without recurrence? No	0.73	0.1536
PT resected without recurrence? Yes	3.07	0.3378

Notes. Relative hazard is for subjects on MGD and WBRT versus subjects on WBRT only. P-value adjusts for multiple testing.

Table 5. Identification of subgroups in split sample analysis: treatment effects in exploratory sample of 0211 trial.

	Factor	Raw	Value A				Value B			
		p-value	Value	n	RH	p-value	Value	n	RH	p-value
france	9.11	0.027	0.0	84	0.32	0.05363	1.0	24	2.94	0.17954
pr.carbo	10.87	0.048	0.0	83	0.39	0.07833	1.0	25	4.26	0.18103
ptstagen	0.21	0.059	3.0	--	1.72	0.42139	5.0	--	0.37	0.07388
neurbl20	12.67	0.088	0.0	98	0.63	0.28314	1.0	10	7.99	0.14457
ltmptbm	0.20	0.098	0.0	49	1.32	0.61953	1.0	59	0.27	0.09389
sympbl12	4.67	0.100	0.0	--	0.47	0.12400	2.0	--	2.18	0.36172
neurbl18	0.22	0.110	5.0	--	0.63	0.29587	5.0	--	0.63	0.29587
basetrla	2.08	0.130	-0.3	--	0.51	0.15601	2.7	--	1.06	0.90295
euroaus	3.86	0.130	0.0	64	0.34	0.11599	1.0	44	1.33	0.61980
metachro	4.25	0.130	0.0	58	0.28	0.11374	1.0	50	1.19	0.73841
prior	3.66	0.140	0.0	67	0.40	0.13240	1.0	41	1.48	0.53320
neurbl21	2.44	0.150	0.0	--	0.57	0.24937	1.0	--	1.40	0.54303
sympbl7	4.71	0.150	0.0	--	0.57	0.19931	0.0	--	0.57	0.19931
prantflg	5.46	0.160	0.0	86	0.54	0.20250	1.0	22	2.97	0.33156
bmlesn	5.17	0.170	0.0	25	0.20	0.14510	1.0	83	1.01	0.98116
sympbl2	1.96	0.170	0.0	--	0.47	0.16551	1.0	--	0.92	0.84827
f.prrad	0.26	0.180	0.0	84	1.02	0.95880	1.0	24	0.27	0.13672
pr.casca	3.23	0.180	0.0	68	0.42	0.15428	1.0	40	1.37	0.61691
karnofsk	1.80	0.190	80.0	--	0.51	0.16114	90.0	--	0.92	0.85675
ltmptdx	0.28	0.190	0.0	51	1.09	0.86437	1.0	57	0.31	0.14175
neversmk	4.81	0.210	0.0	101	0.58	0.23766	1.0	7	2.80	0.37554
prantday	1.46	0.210	0.0	--	0.55	0.21645	44.5	--	0.80	0.62516
pbm.all	4.57	0.220	0.0	100	0.59	0.24367	1.0	8	2.68	0.39687
txgt1.4	3.57	0.220	0.0	90	0.58	0.23659	1.0	18	2.06	0.43119
controll	3.46	0.230	0.0	86	0.57	0.22919	1.0	22	1.98	0.45638
motorlbl	0.38	0.240	5.0	--	0.71	0.45018	5.0	--	0.71	0.45018

Notes. Value A is 0 for binary variables and the first quartile for continuous variables. Value B is 1 for binary variables and the third quartile for continuous variables. Values are specified in the columns labeled value. Sample size is given in columns labeled “n”. RH = relative hazard for MGd plus whole brain radiation therapy (WBRT) versus WBRT. P-value for Value A and Value B adjusts for multiple testing. Factor gives the ratio of variance at Value A and Value B. P-value for factor does not adjust for multiple testing.

Table 6. Distribution of p-values for "uncontrolled" subgroup treatment effects, by sample, in various sample splits.

	P-value ranges					Total
	0-0.01	0.01-0.05	0.05-0.10	0.10-0.25	0.25-1	
Exploratory sample	16	19	11	19	35	100
Validation sample (raw p-values)	14	18	2	1	0	35
Validation sample (adj p-values)	2	9	9	10	4	35

Notes. First row gives distribution of p-values for “uncontrolled” subgroup across 100 sample splits after adjusting for multiple testing. Second and third rows gives distribution of p-values in the validation sample for the 35 splits where uncontrolled group is selected in exploratory sample, i.e., adjusted p-value for uncontrolled subgroup in exploratory sample is less than 0.05. P-values in second row do not adjust for multiple testing. P-values in the third row apply an overly conservative Bonferroni adjustment for multiple testing.

Readers with comments should address them to:

Professor Anup Malani
University of Chicago Law School
1111 East 60th Street
Chicago, IL 60637
amalani@uchicago.edu

Chicago Working Papers in Law and Economics
(Second Series)

For a listing of papers 1–399 please go to Working Papers at <http://www.law.uchicago.edu/Lawecon/index.html>

400. Shyam Baganesh, Foreseeability and Copyright Incentives (April 2008)
401. Cass R. Sunstein and Reid Hastie, Four Failures of Deliberating Groups (April 2008)
402. M. Todd Henderson, Justin Wolfers and Eric Zitzewitz, Predicting Crime (April 2008)
403. Richard A. Epstein, *Bell Atlantic v. Twombly*: How Motions to Dismiss Become (Disguised) Summary Judgments (April 2008)
404. William M. Landes and Richard A. Posner, Rational Judicial Behavior: A Statistical Study (April 2008)
405. Stephen J. Choi, Mitu Gulati, and Eric A. Posner, Which States Have the Best (and Worst) High Courts? (May 2008)
406. Richard H. McAdams and Janice Nadler, Coordinating in the Shadow of the Law: Two Contextualized Tests of the Focal Point Theory of Legal Compliance (May 2008, revised October 2008)
407. Cass R. Sunstein, Two Conceptions of Irreversible Environmental Harm (May 2008)
408. Richard A. Epstein, Public Use in a Post-*Kelo* World (June 2008)
409. Jonathan R. Nash, The Uneasy Case for Transjurisdictional Adjudication (June 2008)
410. Adam B. Cox and Thomas J. Miles, Documenting Discrimination? (June 2008)
411. M. Todd Henderson, Alan D. Jagolinzer, and Karl A. Muller, III, Scierer Disclosure (June 2008)
412. Jonathan R. Nash, Taxes and the Success of Non-Tax Market-Based Environmental Regulatory Regimes (July 2008)
413. Thomas J. Miles and Cass R. Sunstein, Depoliticizing Administrative Law (June 2008)
414. Randal C. Picker, Competition and Privacy in Web 2.0 and the Cloud (June 2008)
415. Omri Ben-Shahar, The Myth of the “Opportunity to Read” in Contract Law (July 2008)
416. Omri Ben-Shahar, A Bargaining Power Theory of Gap-Filling (July 2008)
417. Omri Ben-Shahar, How to Repair Unconscionable Contracts (July 2008)
418. Richard A. Epstein and David A. Hyman, Controlling the Costs of Medical Care: A Dose of Deregulation (July 2008)
419. Eric A. Posner, *Erga Omnes* Norms, Institutionalization, and Constitutionalism in International Law (August 2008)
420. Thomas J. Miles and Eric A. Posner, Which States Enter into Treaties, and Why? (August 2008)
421. Cass R. Sunstein, Trimming (August 2008)
422. Cass R. Sunstein, Second Amendment Minimalism: Heller as Griswold (August 2008)
423. Richard A. Epstein, The Disintegration of Intellectual Property (August 2008)
424. John Bronsteen, Christopher Buccafusco, and Jonathan Masur, Happiness and Punishment (August 2008)
425. Adam B. Cox and Thomas J. Miles, Judicial Ideology and the Transformation of Voting Rights Jurisprudence (August 2008)
426. Daniel Abebe and Jonathan S. Masur, International Agreements and Internal Heterogeneity: The “Two Chinas” Problem (August 2008; updated September 2009)
427. William Birdthistle and M. Todd Henderson, One Hat Too Many? Investment Desegregation in Private Equity (August 2008)
428. Irina D. Manta, Privatizing Trademarks (abstract only) (September 2008)
429. Paul J. Heald, Testing the Over- and Under-Exploitation Hypothesis: Bestselling Musical Compositions (1913–32) and Their Use in Cinema (1968–2007) (September 2008)
430. M. Todd Henderson and Richard A. Epstein, Introduction to “The Going Private Phenomenon: Causes and Implications” (September 2008)
431. Paul Heald, Optimal Remedies for Patent Infringement: A Transactional Model (September 2008)
432. Cass R. Sunstein, Beyond Judicial Minimalism (September 2008)
433. Bernard E. Harcourt, Neoliberal Penalty: The Birth of Natural Order, the Illusion of Free Markets (September 2008)
434. Bernard E. Harcourt, Abolition in the U.S.A. by 2050: On Political Capital and Ordinary Acts of Resistance (September 2008)

435. Robert Cooter and Ariel Porat, Liability for Lapses: First or Second Order Negligence? (October 2008)
436. Ariel Porat, A Comparative Fault in Defense Contract Law (October 2008)
437. Richard H. McAdams, Beyond the Prisoners' Dilemma: Coordination, Game Theory and the Law (October 2008)
438. Dhammika Dharamapala, Nuno Garoupa, and Richard H. McAdams, Belief in a Just World, Blaming the Victim, and Hate Crime Statutes (October 2008)
439. M. Todd Henderson, The Impotence of Delaware's Taxes: A Short Response to Professor Barzusa's *Delaware's Compensation* (October 2008)
440. Richard McAdams and Thomas Ulen, Behavioral Criminal Law and Economics (November 2008)
441. Cass R. Sunstein, Judging National Security post-9/11: An Empirical Investigation (November 2008)
442. Eric A. Posner and Adrian Vermuele, Crisis Governance in the Administrative State: 9/11 and the Financial Meltdown of 2008 (November 2008)
443. Lee Anne Fennell, Adjusting Alienability (November 2008)
444. Nuno Garoupa and Tom Ginsburg, Guarding the Guardians: Judicial Councils and Judicial Independence (November 2008)
445. Richard A. Epstein, The Many Faces of Fault in Contract Law: Or How to Do Economics Right, without Really Trying (December 2008)
446. Cass R. Sunstein and Richard Zeckhauser, Overreaction to Fearsome Risks (December 2008)
447. Gilbert Metcalf and David Weisbach, The Design of a Carbon Tax (January 2009)
448. David A. Weisbach, Responsibility for Climate Change, by the Numbers (January 2009)
449. M. Todd Henderson, Two Visions of Corporate Law (January 2009)
450. Oren Bar-Gill and Omri Ben-Shahar, An Information Theory of Willful Breach (January 2009)
451. Tom Ginsburg, Public Choice and Constitutional Design (January 2009)
452. Richard Epstein, The Case against the Employee Free Choice Act (January 2009)
453. Adam B. Cox, Immigration Law's Organizing Principles (February 2009)
454. Philip J. Cook, Jens Ludwig, and Adam M. Samaha, Gun Control after *Heller*: Threats and Sideshows from a Social Welfare Perspective (February 2009)
455. Lior Jacob Strahilevitz, The Right to Abandon (February 2009)
456. M. Todd Henderson, The Nanny Corporation and the Market for Paternalism (February 2009)
457. Lee Anne Fennell, Commons, Anticommons, Semicommons (February 2009)
458. Richard A. Epstein and M. Todd Henderson, Marking to Market: Can Accounting Rules Shake the Foundations of Capitalism? (April 2009)
459. Eric A. Posner and Luigi Zingales, The Housing Crisis and Bankruptcy Reform: The Prepackaged Chapter 13 Approach (April 2009)
460. Stephen J. Choi, G. Mitu Gulati, and Eric A. Posner, Are Judges Overpaid? A Skeptical Response to the Judicial Salary Debate (April 2009)
461. Adam B. Cox and Eric A. Posner, The Rights of Migrants (April 2009)
462. Randal C. Picker, The Google Book Search Settlement: A New Orphan-Works Monopoly? (April 2009, revised July 2009)
463. Randal C. Picker, The Mediated Book (May 2009)
464. Anupam Chander, Corporate Law's Distributive Design (June 2009)
465. Anupam Chander, Trade 2.0 (June 2009)
466. Lee Epstein, William M. Landes, and Richard A. Posner, Inferring the Winning Party in the Supreme Court from the Pattern of Questioning at Oral Argument (June 2009)
467. Eric A. Posner, Kathryn Spier, and Adrian Vermeule, Divide and Conquer (June 2009)
468. John Bronsteen, Christopher J. Buccafusco, and Jonathan S. Masur, Welfare as Happiness (June 2009)
469. Richard A. Epstein and Amanda M. Rose, The Regulation of Sovereign Wealth Funds: The Virtues of Going Slow (June 2009)
470. Douglas G. Baird and Robert K. Rasmussen, Anti-Bankruptcy (June 2009)
471. Bernard E. Harcourt, Alon Harel, Ken Levy, Michael M. O'Hear, and Alice Ristroph, Randomization in Criminal Justice: A Criminal Law Conversation (June 2009)
472. Bernard E. Harcourt, Neoliberal Penalty: A Brief Genealogy (June 2009)
473. Lee Anne Fennell, Willpower and Legal Policy (June 2009)

474. Richard A. Epstein, How to Undermine Tax Increment Financing: The Lessons of ProLogis v. City of Chicago (June 2009)
475. Randal C. Picker, Online Advertising, Identity and Privacy (June 2009)
476. M. Todd Henderson, Credit Derivatives Are Not “Insurance” (July 2009)
477. Lee Anne Fennell and Julie Roin, Controlling Residential Stakes (July 2009)
478. Douglas G. Baird, The Holmesian Bad Man’s First Critic (August 2009)
479. Douglas G. Baird, The Bankruptcy Exchange (August 2009)
480. Jonathan Masur and Eric A. Posner, Against Feasibility Analysis (August 2009)
481. Lee Anne Fennell, The Unbounded Home, Property Values beyond Property Lines (August 2009)
482. Bernard E. Harcourt, Henry Louis Gates and Racial Profiling: What’s the Problem? (September 2009)
483. Stephen J. Choi, Mitu Gulati, Mirya Holman, and Eric A. Posner, Judging Women (September 2009)
484. Omri Ben-Shahar, One-Way Contracts: Consumer Protection without Law (October 2009)
485. Ariel Porat, Expanding Liability for Negligence *Per Se* (October 2009)
486. Ariel Porat and Alex Stein, Liability for Future Harm (October 2009)
487. Anup Malani and Ramanan Laxminrayan, Incentives for Surveillance of Infectious Disease Outbreaks (October 2009)
488. Anup Malani, Oliver Bembom and Mark van der Laan, Accounting for Differences among Patients in the FDA Approval Process (October 2009)