

(5)世界的に見た日本の臨床試験のレベルとその推移

司会：谷田憲俊（兵庫医科大学第4内科）

編集部：このセッションは、昼休みに読んでおいてくださいということで、英文論文の一部（2論文）と日本語論文を参加者に配付しました。午後の、論文評価のハダッドスケールで使用するためです。柳氏の話に先立って、浜から論文抽出までの経緯の紹介をしました。

浜：日本の臨床試験を世界的なレベル、一番簡単な基準で一度ざっと見てみるということをしよつと考えたのですが、昨年のオーストラリアのアデレードで開かれたコクラン・コロキウムに参加したときに、カナダのMoherさんという生物統計学者で疫学者でもある方にお会いして相談しましたら、1980年、1990年、1995年の臨床試験論文を50ずつくらいを無作為抽出して、ハダッドスケールという比較的確立したものを使って評価してみたらどうか

とサジェスチョンを受けました。たまたま、92年版二重盲検文献資料集という本が手元にありまして、二重盲検の論文数が6件以上あった11雑誌を抜き出しました。これらの雑誌に載っていた80年、90年、95年の論文を取り出したわけです。「比較試験」と銘打っているものは取り出しました。第3相の臨床試験と銘打っていて比較試験になっていないというものもありました。これはいったい何だと思いました。初めて体験いたしました。第3相の臨床試験は当然、比較試験で確認すべきと思いますが、そういうことで第3相となっておれば比較試験でなくても含めました。

その後TIPのメンバーでそれらの雑誌から論文をハンドサーチし、乱数表を使って、50論文を抽出しました。それを4人で、独自に、ハダッドスケールで判定しました。その結果をこれからお話します。

柳 元和 ———— 内科医師

ハダッドスケール

元々は痛みの臨床試験について、本当に意味のあるスクリーニングができるかどうかということで開発されたスケールです。測定項目をたくさん作って、その中で臨床経験のある医師の意見も突き合わせて、意味のある項目を選抜していく。最後まで残った項目で、トレーニングを受けていない人でも、このスケールを使えばある程度はきちんとした評価をできる、そういう作業を経てこのスケールが構成されています。

そのスケールは痛みの研究についてはきちんとしていますが、一般の臨床試験について使え

るかどうかは確認されていません。しかしながら、われわれ日本の文献を概括するというのを第一の目標にして、これを使って検索することにします。

採点基準

第一には、「無作為」あるいは「無作為化」と書いてあるかどうか。「二重盲検」あるいはこの概念に相当する言葉が書いてあるかどうか。「解析除外例」「脱落例」について詳細に1例も漏らさず記載されているかどうか。それぞれ、「イエス」であれば、「1点」を与える。

付加事項がありまして、質問1.無作為であるということに対して何の記述もなければ、1点を与えておこう。しかし、内容を読んでみて、これはランダム化ではないと考えれば「-1点」として、±ゼロになる。二重盲検についても同様に、ブラシーボを使っていないとかダブルダミーにしていない、活性物質が2つある場合にはダミーも2つ必要ですが、そういうことがきちんと記載されていない場合には、不適切と判定できる。つまり、書いていない場合には判定のしようがないが、書いてあって不適切な場合は「-1点」。逆に、手順をきちんとしている場合には「+1点」加算で、合計2点にする。そういうことで、大まかに論文の点数にバラツキが出てきます。

3点を境にして、4点、5点は、かなり手順としてきちんとしている論文という判定を下せます。1点や2点の論文は問題がある。論文に記載がないか、方法論として認められないようなものだと判定できます。

判定の具体例

具体的には、ランダム化、どちらの組に入るかが事前にまったく予測できないことを保証する方法については「適切」と認めましょう、しかしそれ以外の方法、誕生日、入院日、交互に割り付ける、こういう方法は適切と言えない。次に来る患者さんはこっちだな、と予測できてはいけない、これは最も基本的なことです。二重盲検についても、飲んでみたらブラシーボと分かるのでは話にならない。きちんと書いているけれども、実は嘘を書いていた、ということは見破ることは出来ませんが、一応文面については信用せざるを得ない。

脱落についても理由が分からないのは、厳密には「脱落」とは言わない。後追いの調査怠慢である。具体的には、一人ひとりの患者さんを追いかけて行って、どういうときにどうなったとききちんと記載していないと、きちんとしたトライアルとは言えない、というのがコクランなどではそういう考え方です。もしも「来院せず」

とあった場合にどうするかは（今回の4人の採点者間で）議論できていないのですが。私は、「来院せず」とあるだけでは「0点」だと思うのですが。何の情報も得られないわけで、その人に対してどういうフォローアップであったかが分からない。

論文の読み方——エパルレスタット論文を例に

無作為化の論文を読むときに、通常、われわれはサマリーを読みます。メドラインでサマリーだけなら手に入ります。これを読んで、だいたいの当たりを付けます。みなさんにお配りしている資料の英文論文では、サマリーに、「Duble blind：ダブルブラインド」という言葉が出てきます。しかし「Randomization：ランダム化」はどこにもない。どうもよくわからないな、となりますと、中身を読んでいかなくてはならない。中身を読むとなると、小見出しが大事です。「Introduction：イントロダクション」には書いていないです。次の小見出しにPatients and Methods方法論が出てくる。どういう人を対象としているかがわかる。ところがここを読んでRandomizationのことは書いていない。

この論文は、エパルレスタットが日本で承認されるに際して参考にした日本語文献を英語に置き換えたのです。実は、日本語の論文には、「無作為割り付けをした」と書いてあります。何症例かをひと組にして、それをランダム化して、各施設に配りましたという表現です。これは従来の論文でランダム化というときに、コクランでいうところのランダム化には当たらない。おそらくこの英文を作成するときに、Randomizationとして送ったのではないかと、しかし、削られたのではないかと私は思います。この論文は、アルドース還元酵素剤のレビューでは、採用されませんでした。おそらく査読した人が、Randomizationと認めないと言ったのではないかと想像しております。

論文の読み方——トルレスタット

もう一つの英文論文（資料参照）はトルレスタット

タットです。サマリーの中に、Randomizationとあります。サマリーを読んだだけで分かります。あとは中身を読んで確認をしていく。小見出しの「Drug administration」にも randomisedという言葉があります。しかしこれは、ずっと全体を読んだ限りでは問題があると私は思います。Randomizationしているからといって、すべて良いとは限らない(このことの議論は今回は省きますが)。結論だけを申しますと、この論文はトルレスタットをかなり有望だと言っているのですが、昨日もお話しましたように(H.糖尿病用剤)、トルレスタットは開発が中止されました。このような論文が一つあるからといって、信用するわけにはいかない。

ハダッドスケール評価例—その1

アトピー性皮膚炎に対する臨床効果という論文は、サマリーは英文です。これを読んだだけでは分かりません。本文を読みますと、英語論文と同じで、小見出しは、「はじめに」「対象および方法」...とあり、最後に「全般改善度」というのがあります。「試験方法」というのがあります。これには、プラシーボを使ったかとか、ランダム化されたかとか、まったく書かれていません。これで第3相長期試験です。

認可にかかわる非常に重要な、しかもアトピーですから長期に服用しなくてはならないだろう、そういうものでランダム化されたものの記載がない、プラシーボの記載がない。どういふことかとずっと読んでいきますと、20mg、30mg、40mg、症例の構成という記載があります。やはり比較試験のようです。なぜ比較出来るのか理解できなかったのですが、主治医が投与量を決めるらしいのです。方法論に「治験薬は症状によって適宜増減してもよい」と書いてあります。投与量は主治医が決めて、その各ミリ数間で比較するという、比較試験です。この評価を4人がどうだったかは後で紹介します。

ハダッドスケール評価例—その2

塩酸イブサピロン、不安神経症に対する臨床

評価の論文です。これは前期第2相、用量設定などを決めていく試験です。これも、小見出しは、「はじめに」「対象および方法」...とあって、どういう対象を選ぶか、インフォームド・コンセント、薬剤、スケジュールなどが書いてあります。形式は整っているように思います。試験薬剤を見てみますと、「外観上、識別不能な2.5mg錠、5mg錠、プラシーボの3剤を使った」と記載があります。識別不能とはっきり書いてありますから、一応文章上は信用せざるを得ない。投与量および投与方法については、「シングル・ブラインド」と書いてあります。ダブル・ブラインドではないことがわかります。

論文評価の第一関門—せめて形式を

つまり、形式だけでだいたいその試験の持つ臨床的な意味の重みが推定できる、ということです。しかし、これは第一段階です。論文の最終的な質とは関係ありません。たとえ形式を整えていても質的に悪い論文はあります。しかし、逆はないのではないのでしょうか。形式が整っていないのに、良い質の論文が出るだろうか。そこにわれわれは着目したのです。ですからこれは論文評価の第一関門です。

こういうチェック項目を作って論文を読んでいけば、書いているかどうかは分かる、つまり医学に素人の方々でも分かる、チェックできる。例えば薬を病院でもらった。その薬についての論文を取り寄せた。全部を読むのは大変ですが、少なくとも、ランダム化されているか、プラシーボを使っているかはチェックできる。もしも使っていないとなると、この薬はちょっとおかしいのじゃないか、と疑問を持つきっかけになるだろう、というのがわれわれの狙いです。

慢性心不全におけるエナラプリルマレイト

これは2つの試験を一つの論文にまとめている。1回投与における運動能の改善と、長期投与における「全般改善度」が出ています。プラシーボを使い、識別不能としています。問題は、9

症例分を1組とする割り付けを行ったとあります。これが本当にランダムイズド・コントロール・スタディかどうかは議論の対象になると思います。

患者背景に、NYHA分類で2度、3度の患者ばかりです。4度つまり重症の人はいないということです。

結論は、試験物質とプラシーボとで差がなかったんです。しかし、エナラプリルは国際的に認められている良い薬です。国際的に認められている薬剤が、「良い」という結論を出せなかった論文です。これは日本の文献で、良い薬を良いと言えなかった例です。有意差が出せなかったんです。この論文を読むと、日本の臨床家は、国際的に認められているエナラプリルに対して、あんまり効いていないと思うかもしれない。そういうおかしなことになる例として紹介しました。

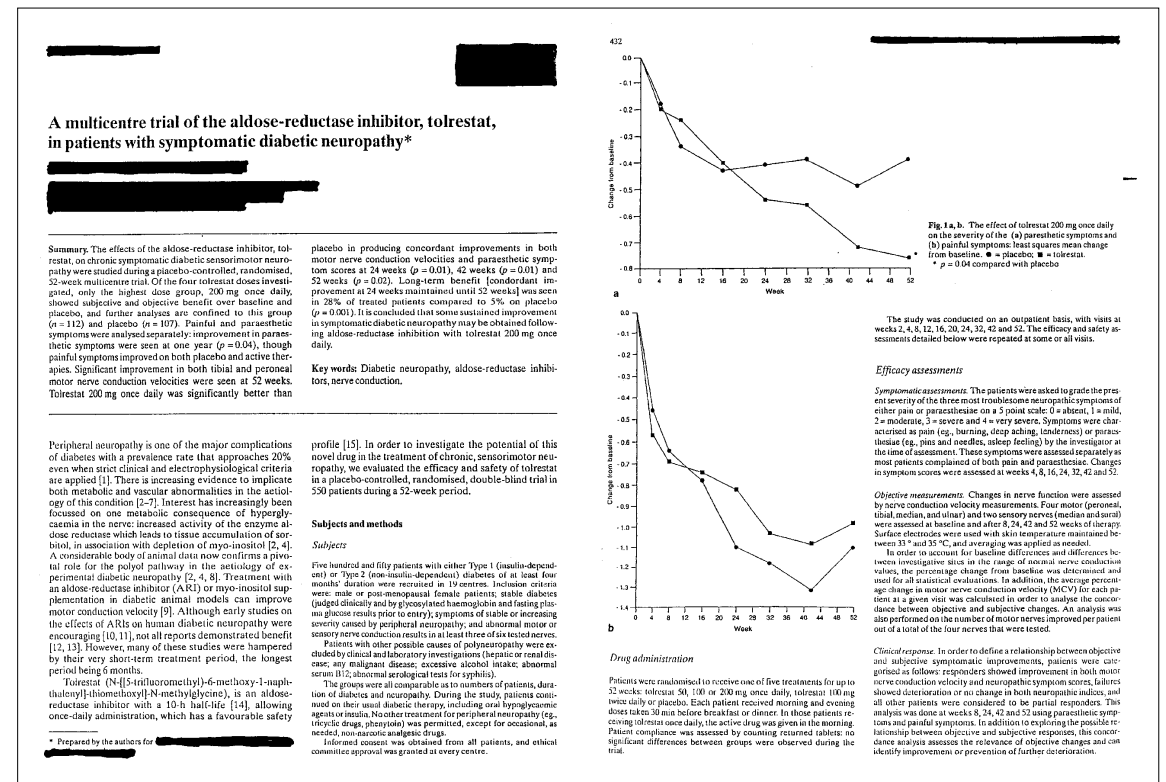
4人の評価結果

4人とも評価が一致していたのが、1980年は17

件、1990年は15件、1995年は24件です。95年になると論文の半分くらいは一致している。4人のうち1人は全く医学の素人です。それでも4人一致するのが24件も出てくるというのは、スケールが単純であることで比較的判定しやすい。1点や2点の論文は、このスケールでは判定しがたい要素があるということです。このスケールには限界がある。その分析はまだできていません。しかし、何のトレーニングも受けていない4人がこれだけ一致しているというのは、大まかなところでは有望ではないか、スケールとしてかなり使えるのではないかと思います。

みなさんに資料をお配りした論文の45番(エパルレスタット)は、1点、0点、1点、1点でした。ひどいです。97番(トルレスタット)は3点、1点、1点、1点です。一人は評価が甘いですね。122番(塩酸イブサピロン)は、3点、2点、2点、2点。きょうご紹介した3論文は、点数が低いところで一致している。

5点は不思議と4人が一致しますね。1980年に



資料

5点が一つありました。平均値はこういう採点に意味があるかどうかは別として、目安として報告します。1980年が4人の平均の平均が2点くらい。1990年は落ちまして、1.5点くらいです。1995年が挽回して1.7点くらいに回復してきました。1990年というのは、臨床試験として(質が)落ち込む時期だったのではないかと想像しております。昔のほうむしろ良かった。90年代に入って、臨床試験の質が非常に低下している可能性がある。それがいろいろ批判を受けて、今、回復期にあるということではないでしょうか。

浜: 柳さんの推測は、あながち間違いではないのです。薬事法が改正されたのが1979年です。それ以前から行われていた臨床試験が発表されるのが1980年で、大きな薬害が表沙汰にならないできた1990年以降に承認される薬剤がかなりルーズになってきたことと、多少関係するかと思えます。

柳: 1990年代に新薬として認可されたものは危ない、と思ったほうがいい。

浜: 1995年頃に承認された薬剤というのは、1990年頃に臨床試験の報告がなされたものが多いですから。承認すべきでなかったような問題の新薬は94年、95年よりは96年のほうが少なかったです。多少、今後はましになるのかなと考えます。

谷田: 最後のプレゼンテーションは非常に分かりやすいものでしたから、みなさん、消化不良にはならず済んだかと思えますが、却って暗澹たる思いになるかもしれません。製薬会社が一方的に責められているように思われるかもしれませんが、実際に責められるのは日本の医者、研究者です。いかに今まで、いい加減なことしてきたかを認識していただくものだったと思います。医療関係者だけでなく他の方々も含めて、自らのレベルを上げていかなくてはいけないということをこのデータは示していると思います。

■ 質疑 / 討論 ■

近藤: (医師、放射線科) 質問です。イブサピロンの不安神経症に対する論文は前期第2相ですから、もともと一種の増量試験だと思うのです。ですからダブルブラインドや無作為化を要求するのは無理なのではないかと思いましたが、どうなのでしょう。僕の知識が足りないのかもしれませんが。

柳: 完全な用量設定試験であれば、ダブルブラインドにする必要はむしろありません。だからプラシーボを使う必要もないと思うのです。

近藤: プラシーボというのは、錠剤を...

柳: こういうデザインで行われていること自体がおかしいのです。完全な用量設定試験であればプラシーボと比較するということは...

近藤: はい、わかりました。第2点は、著者名や施設名をブラインド(編集部: 参加者に資料として配った文献の著者や施設名を黒塗りしたこと。資料参照)にしたのはどういう理由からですか。

柳: 施設とか著者の名前によって、採点者がこれは良い論文かどうかと先入観で見ないようにです。(会場爆笑)

近藤: さっきの浜さんの話で少し奇異に感じたのは、評価したのがどの薬なのか、名前が出て来ないことです。

浜: 最後のほうでちらっとお見せしました。

近藤: はい、少しね。どこかで発表されておれば別ですが。

浜: いえ、まだです。

近藤: 匿名の批判は批判足りえないはずですが、批判する相手をきちんと出していかないのも批判にならないのではないかと思うのです。(柳: これは批判のためではなく評価の仕方のセッションですが。) ええ、そのことは分かります。ただ、評価するにも匿名の文献

というのどうかなと思ったものですから。データがブラインドである理由は先ほどの説明でわかりましたが、最初にそういう説明がほしかったと思います。

柳: はい。雑誌名とか著者の名前で論文に対する評価が、かなり引きずられるという研究があるんです。ですから敢えてブラインドにさせていただきました。97番の論文について補足するのを忘れました。「脱落」「中止」の内訳をちゃんと書いています。しかし、53症例中31例が脱落とはどういうことか。(会場笑い) けしからんじゃないか、というのが私の評価です。

浜: もう一度、JIP/TIPが評価した95年~96年承認の新薬リストをお見せします。この中にはヒスマナルもあります。これは抗アレルギー剤でトリルダンと非常によく似た薬剤です。これ自身ももちろん(副作用として)QT延長から不整脈を起こすんですが、抗アレルギー剤としてどれか一つ残すとすれば、まあこれかな、と。しかし非常に注意して使わなくてはいけないということで、5ではなくて4の評価にしました。95年には、結構、6つまり駄目だというのがあります。

今日はみなさんに一応お見せするために間に合わせようと非常に粗っぽい段階のものしか準備できませんでした。実際の評価しようとする、非常なエネルギーが要ります。もしもみなさんの中に一緒に評価しようという方がおられましたら、是非手伝っていただきたいと思います。

佐藤: (東京医科歯科大学) 非常に重要なことですので発言を。97番と45番の論文は、そもそもコントロールを置いた比較試験ではない論文なのです。それをもってきて批判するのは全くナンセンスでして...

谷田: ちょっと待ってください。これは批判ではなくてトレーニング練習問題でして、評価項目を見ていただいたら分かると思うの

ですが、トレーニングなのです。

佐藤：だから、トレーニングであるにしてもあまりにも不適切な教材を取り上げたと思います。私の専門は主に市販後調査です。臨床試験のことは必ずしも専門ではない。先ほど近藤さんもおっしゃいましたが、そういう私にさえ、教材として不適切ということがわかるというのは非常に残念な発表だと思います。

例えば97番について。そもそもコントロールを置いた比較試験ではなくて、3群に分けていますが、あくまでも用量関係をみるためにプラシーボを置いているわけで。

谷田：これはコントロール・スタディではない、ということが分かればいいのですが。その背景はまた別のことでして。

佐藤：いや、論文の評価ということではなくて、本来そうではないスタディに対してダブルブラインドかどうかと判断させるのはどうかと思うのです。つまり二重盲検なり比較試験でない論文を取り上げて...

谷田：ですから、この論文を読んで二重盲検でないということが分かっていたらいいのです。それから先は別の機会に話していただけますか。いろいろ消化不良の点はありますが、時間が迫っておりますので。(会場で拳手あり)あ、まだ(会場笑い)はい。

津谷：122番の論文ですが、実施方法のところで「割り付けは4症例分を1組とし、割り付け表はコントローラーが保管した」。この「割り付けは4症例分を1組」というのが何のことか分からないとおっしゃっていましたが、これはいわゆる「ブロック・ランダム化」というものです。割り付けの基本中の基本です。私の講義を受けた人は、学生も消費者グループの人もちちゃんと知っています。つまり各施設に、この場合には2例、2例が行くようにという、施設をブロックとして割り付けたということです。

柳：それ、世界的な方法と考えられていますか。

津谷：はい、世界的な方法です。

柳：何十という施設に対して。

津谷：ええ、そうです。

柳：この点についてはまた議論したいと思います。

谷田：TIPの評価というのは、「物質」から「薬」と名前を変えたものが、果して本当に「薬」として承認してよかったのかどうか、またその「薬」が実際に臨床で本当に役に立つかどうかを判断するものです。

新しいGCPが今年4月から導入されて、今まで指摘されていたような問題は起こらないはずなのですが、日本の医療状況や研究者、医師、医薬品業界を含めて、果して新GCPに対応できるのだろうか考えると、疑問の部分があると思うのです。そういう部分は、市民や患者さんが監視していかなければならないことだと思います。情報の公開に関して、次のセッション(P.460)で話がありますのでみなさんに十分にお考えいただきたいと思います。

本当に申し訳ないのですが、時間がないのでこのセッションはこのへんで終わりにさせていただきます。

編集部：ブロックランダム化は世界的に行われている方法であるが、1998年11月に厚生省医薬安全局から示された新しい「統計原則」ではブロックを十分長くし、2種以上の長さのブロックを用意し、ブロックの長さが知られないようにするなど、中央で管理すべき内容について詳しく規定している。2例、2例、合計4例を各施設に送り付けた後、中央は関知しないという従来の方法は明瞭に否定されている。

TIP Vol.14, No.2, P.15-19, 1999 参照

第7章

最近の課題、医療事故 (K)

[1] 予防接種(ワクチン)

藤井俊介
山本英彦

[2] 陣痛促進剤(子宮収縮剤)

勝村久司
石川寛俊

[3] アトピーとステロイド剤

深谷元継
玉置昭治
林 敬次
津田敏秀
水間典昭

[4] 最近のトピック

フェノテロールと
喘息患者の死亡
土居 悟
浜 六郎
福本真理子