

266400960A

厚生労働科学研究研究費補助金

医療技術評価総合研究事業

医師国家試験コンピュータ化に関する研究

平成16年度 総括・分担研究報告書

主任研究者 細田 瑛一

平成17(2005)年4月

## 目 次

### I. 総括研究報告

医師国家試験コンピュータ化に関する研究	-----	1
細田 瑛一		
高林克日己		

### II. 分担研究報告

1. コンピュータ試験問題作成と運用の研究	-----	15
高林克日己		
2. 試験問題のトライアルの実施と検討	-----	23
福島 統		
3. 試験問題のトライアルの実施と検討	-----	25
吉岡俊正		
4. 試験問題のトライアルの実施と検討	-----	26
江口光興		
5. 試験問題のトライアルの実施と検討	-----	27
杉山幸比古		
6. 試験問題のトライアルの実施と検討	-----	29
山内俊雄		
7. 試験問題のトライアルの実施と検討	-----	31
原田研介		

III. 研究成果の刊行に関する一覧表	-----	33
---------------------	-------	----

IV. 研究成果の刊行物・別刷	-----	34
-----------------	-------	----

厚生労働科学研究費補助金（医療技術評価総合研究事業）

総括研究報告書

医師国家試験コンピュータ化に関する研究

主任研究者 細田 瑳一 財団法人日本心臓血圧研究振興会附属榊原記念病院病院長  
高林克日己 千葉大学医学部附属病院企画情報部教授

研究要旨

【目的】国家試験にコンピュータを利用した試験を導入する蓋然性を検討する。【方法】MCQ問題100題と患者管理問題（PMP）5題、及びPMP問題に対応する内容のMCQ問題10題を作成し、6つの大学の医学生5、6年生を2群に分けて従来のマークシート法（MAR）とコンピュータ試験（COM）で半分ずつを行い、全ての結果を総括し解析することでCOMとMARの相関と利点欠点を検討する。【成績】COMとMARのMCQに関する解答の結果は近似しており、試験運用上でも大きな問題はなかった。PMPのCOMの結果とMCQ50題の正答率との間には評価項目及び解析可能な内容が異なるので強い相関はみられなかった。またPMPに対応するMCQ問題とも相関が弱く、PMPでの評価判定は知識のファクターだけではないことが推測された。またPMP問題ではとくに能力の低い学生の識別が明瞭であった。しかし多数の受験生の同時受験などでの運用上の問題が解決されていないことと、学生自身はCOMに賛成する者より反対意見を述べる者がMCQで82:381、PMPで116:402と圧倒的に多かった。決定的に有用な事由がないと早急にコンピュータ試験を導入することには困難があると考えられるが、試験自身には問題のないことが示されたことから、CATの導入が可能になるなどさまざまな利点を考慮して、欠点を克服した段階で近い将来導入されることは十分考えられる。【結論】試験方法としては国家試験にコンピュータを導入し、従来のマークシートに変わり共用試験のCBTのように試験ができること、PMPのようにペーパーテストで検定できない能力の判定ができる可能性を示した。

分担研究者 高林克日己

千葉大学医学部附属病院  
医療情報部教授

吉岡俊正

東京女子医科大学  
医学教育学 教授

福井次矢

聖路加国際病院院長

江口光興

獨協医科大学  
小児科学教授

杉山幸比古

自治医科大学  
呼吸器内科学教授

山内俊雄

埼玉医科大学  
精神医学教授

福島 統

東京慈恵会医科大学  
解剖学 教授

原田研介

日本大学医学部  
小児科学教授

分担協力者 椎橋実智男

埼玉医科大学  
医療情報施設 助教授

宮木浩行

三菱電機インフォメーションシステムズ株式会社

## A. 研究目的

コンピュータの MCQ 方式と従来のマークシート試験の間にどのような差があるかの検定と、PMP 方式の問題と従来のマークシートによる MCQ 問題との間にどの程度の相関があるかを検定し、運用についての問題点と受験生の対応、利点欠点を示すと共に検討し、コンピュータ試験実施の蓋然性を明らかにすることを目的とする。

## B. 研究方法

### 1 問題作成

昨年度までに作成した CBT に準じたコンピュータでの multiple choice question (MCQ) のツールを用いて、各大学が独自に作成した 5、6 年生の内科系、外科系問題の中から同一の分野でレベルの近似した問題を集めてそれぞれ対の問題になるような 100 題を抽出し、各 A、B 問題 50 題とし、それぞれの問題をコンピュータ用とペーパー試験用に使えるよう作成した。PMP は昨年までに作成した PMP のツールを利用してコンピュータ用に 5 題を高林が作成し、このうち 4 題については各 PMP 問題に対応するように PMP 1 題あたり 5 問の MCQ 問題をペーパー試験用に作成した。

### 2 試験方法

今回は 6 大学において主に 5 学年の学生を対象に試験を行った。試験はコンピュータによる MCQ 試験 50 題 (75 分)、PMP 試験 3 題 (45 分) と、ペーパー試験の MCQ 60 題 (75 分) である。各施設ではグループを対の 2 群 (グループ I、グループ II) に分け、グループ I はコンピュータで MCQ の A 問題の試験、PMP では 3 問を受験し、ペーパー試験では MCQ の B 問題、および PMP で受けなかった 2 問に対応するペーパー試験 10 題を受験した。グループ II はこの逆に MCQ の B 問題の試験、PMP ではグループ I とは別の 2 問と共通の 1 問を受け、ペーパー試験では A 問題、および PMP で受けなかった 2 問に対応するペーパー試験 10 題を受験した。この間 A、B

の受験生は別室でコンピュータ、あるいはペーパー試験を受けるか、あるいは同時に試験を行い、両グループの受験生が試験内容を相互に話し合える時間はないようにした。

### 3 総受験者数

慈恵医大、日大、東京女子医大、埼玉医大、自治医大、獨協医大の 5 年生、(2 大学では 6 年生を含む) を対象に、男性 264 名、女性 216 名、計 580 名を同一施設、同一学年の学生を 2 群に分けて試験を行った。

(倫理面への配慮)

大学名、学生の氏名は用いず、コード番号で表示した。

## C. 研究結果

### 1 MCQ における検討

MCQ 試験による得点を表 1 および図 1 にまとめた。平均点では 6 年生に行なった C 大学と、F 大学の 6 年生が 60 点台であったのに対し、もっとも低い大学は 30 点台と大学群間で乖離した数字となった。コンピュータ試験 (COM) とペーパーテスト (MAR) の結果において、2 群間には強い相関がみられ、平均点で COM の方が 0.24 低い結果が得られた。男性でも女性でもともに同様の結果であり、各大学でみても COM の平均点の方が低かったが、6 年生に行なった F 大学の結果のみが、COM の平均点の方が高かった (ここでは後に述べる 4 問の問題を除外した数値で示してある)。

図 1 をみるとコンピュータ試験 (COM) に比してペーパー試験 (MAR) でとくに得点の低い群と高い群が少数であるが存在する。 $r=0.70$  でほぼよい相関といえる。男女別にみると、図でみると男性の方がより高得点にみえるが、これは 5 年生に女子学生だけの大学が含まれていたためであると考えられる。

次に得点分布を偏差値で比較すると、図 4 に示すように、A 問題、B 問題ともにほとんど二つのグラフは重なる (CBT: コンピュータ試験、MAR: ペーパー

一試験)。

個々の問題についての得点差について検討したのが図5である。ほとんどの問題でコンピュータとペーパーテストの間に正答率に差はないといえる。このうちA8、A50、B27、B44は紙とコンピュータ問題で問題・解答内容に違いがあり除外することにした。A7は選択問題X2であった。したがって他のデータ、図表ではこれらを除外した結果を示してある。B10のみは両者の正答率に大きな差がでているが、その意味を説明できなかった。4問を除外した結果、平均点±標準偏差はコンピュータテスト  $50.88 \pm 15.01$  に対してペーパーテストは  $51.12 \pm 14.50$  となり、平均点ではほぼ同じになった。

図6は個々の得点を2群に分けて比較したものである。平均点、最高点、最低点などはほぼ同じであり、両群間に有意差はない。

## 2. PMP試験(PMP)と従来の紙試験(MAR)との比較

PMPの問題5題のうち3題を施行しているので、全員の3題の平均得点と、全員のそれに対応するMCQ10問の正答率を示す(表2)。

### 1) PMP試験と対応するマークシート試験との比較

3題のPMP問題の得点とMCQの得点との間には強い相関はなかった(相関係数  $r=0.439$ )。男女比で見ると女性にやや相関がみられた(図7)。これに対してMCQの得点とPMPに対応するペーパーテストの10問の得点で比較するとこれでは弱い正の相関が見られる。いわゆるペーパーテストでは全体のMCQと同様に知識を問う試験問題になっていることから、PMPと異なった相関を示すものと考えられる(図8)。またPMPのコンピュータ問題とペーパー試験の直接の対比をみると、この散布図(図9)に示すように、両者間には相関がなかった。

### 2) PMP問題での検討

問題のテーマごとにみるとそれぞれ別の傾向がある。アニサキスでは解答項目が3項目で得点もその3段階に分

かれる。特に中間層は少なく、全然できないか、全て答えたか、どちらかの者が多い。一方喘息は問題の内容と解答枝の項目(治療法)が多種であり、得点分布もこれによって多岐に分かれている。肺梗塞も治療法が多くあるために得点分布も広がるが、途中で大きなギャップがある。イレウスも能力のあるものとなないものの差が明確になっており、かなり能力のあるグループを入れると三群に分かれる(図10、11、12、13)。

## 3. 運用上の問題

途中で作動しないことが埼玉医大で発生した。またインストールにおいて、MCQで82:381 PMPで116:402とLAN上のクライアントに配信する方法では、各施設の設定に影響されて多少の問題が起こった。画像の情報がとくに階段教室などでは周囲の学生に見えてしまい、これがヒントになることが考えられた(図14)。

## D. 考察

今回用いた問題が問題として最善のものだけとは言えないので、明確な結論を出すことはできないが、これまで行われていなかった紙の試験とコンピュータの試験を直接客観的に比較する良い機会であった。

MCQ問題におけるCOMとMARの間には大きな差は見られなかった。個々の問題について、正答率に特徴的な差はなかった。形式によってコンピュータとの比較において差異はなかった。男性の方がややCBTで得点が高い傾向が伺えた。アンケートでは画像などが見やすいとの意見があったのに対して、メモをとりたい、目が疲れるなどの意見が出ており、操作性については大きな問題として上げる意見はなかったが、直接の両者を比較した意見ではペーパー試験の方がよいという意見の方が多かった。紙試験に慣れている受験生の意見として当然のことと思われる。

PMPでの得点と全体のMCQの得点と比較することで、PMPの独自性と蓋然

性を検討した。正の相関は得られたが、とくに PMP では低い得点として存在する群が認められ、より明確に能力の低い学生を検出するものか、統合的意思決定能力に欠けるもの、あるいはコンピュータを苦手とするものを検出しているものと思われた。一方で MCQ の結果とある程度の正の相関がみられることから、全く異なった能力を評価しているわけではないと考える。

PMP 問題と対応する MCQ の問題は必ずしも内容が一致しておらず、相互の比較をすることはできないと考えられた。ひとつの PMP の中で治療や診断の解答数の多いほど、得点分布がひろがり、喘息などをみても、かなり均一に広がっている。これに対して有効回答が 3 つしかないアニサキスはとくに極端な得点分布になっている。また全く回答できない例も多く見られ、ペーパー試験に比べ、得点の低いものの検出力に優れているともいえる。一方で診断名、治療数など解答数により得点分布が大きく異なることから、得点で評価するためにはある程度の解答要求数があったほうがよいのかもしれない。しかし逆に何を持って評価するのかがあいまいになることを考えると、必ずしも細かいほうがよいのではなく、アニサキスの問題のように 3 つくらいのカテゴリーに分けるほうが正確な指標であって、これが PMP の限界なのかもしれない。問診の中で適切なものを選んだかどうかの評価まで加えることでより正しい評価を示すものであるといえるが、従来そこまで分析したシステムはなかったし、問題を作るのが大変であろう。しかしせつかくの総合能力の試験であるとするなら問診、所見、検査、診断、治療の総合点で評価すべきものかもしれない。問診なども、たとえば決定的なポイントになるものだけにつけるなどの工夫をすればよいかもしれない。

アンケート結果では MCQ よりも、PMP の方が好ましい結果も得られているが、操作性の問題で批判がある。特に検索方式であるが、これらは改善の余地があると考えられる。

運用上の問題としてインストールの方法のトラブルであったと思うが、これだけの数を正確に同時に動かすことが求められるとなると、試験のリスクが大きい。端末の故障が起こる確率から 100 台あたり数台以上の予備機の準備は必要である。

PMP は画像の含まれる問題ではこれが周囲に見える状況になると大きなヒントを与えることになりかねない。このためには隔離した環境、あるいは広いスペースが必要になる。

#### E. 結論

MCQ 問題をコンピュータ試験で行なうことに大きな問題はなく、またその得点は紙のテストの結果とほぼ一致した。しかしながらメモ機能などを希望する声があり、また同時に多数の受験者を扱うことの問題は残っている。PMP 問題の結果は MCQ の結果と強い相関を示さず、PMP では統合能力、意思決定能力など別の能力を見ていると考えられた。またとくに能力の低いものを明瞭に判別しえた。

現在ペーパーテストをコンピュータ化することの緊急性はないが、今後多くの試験がコンピュータ化していくことは容易に想像される。メモ機能などいくつかの技術的な改良とともに、大量の問題作成方法、および同時に実施する方式をとるのか、別の日時で試験を行う方式をとるのかの検討が必要のところきている。これにより準備は大きく異なる。

あるいは医師免許の更新制が開始されるのであれば、このようなところから開始するという方法もあるかもしれない。

#### F. 健康危険情報 なし

G. 研究発表

1. 論文発表

なし

2. 学会発表

高林克日己 細田瑛一 福島統ほか  
医師国家試験へのコンピュータ試験の  
導入の検討 第37回日本医学教育学  
会大会（予定）

H. 知的財産権の出願・登録状況（予  
定を含む。）

なし

(表1) MCQ問題のコンピュータ試験とペーパーテストの結果

コンピュータ試験		ペーパー試験	
(COM)		(MAR)	
全体		全体	
被験者数	504	被験者数	498
平均点	50.88	平均点	51.12
標準偏差	15.01	標準偏差	14.50
最高点	92	最高点	90
最低点	17	最低点	18
A大			
被験者数	23	被験者数	23
平均点	44.20	平均点	44.75
標準偏差	9.01	標準偏差	7.86
最高点	60	最高点	62
最低点	27	最低点	30
B大			
被験者数	112	被験者数	112
平均点	42.06	平均点	42.52
標準偏差	11.58	標準偏差	11.48
最高点	69	最高点	78
最低点	17	最低点	22
C大			
被験者数	104	被験者数	104
平均点	61.78	平均点	62.23
標準偏差	12.42	標準偏差	12.45
最高点	90	最高点	90
最低点	21	最低点	24
D大			
被験者数	89	被験者数	86
平均点	54.99	平均点	54.30
標準偏差	11.19	標準偏差	9.64
最高点	79	最高点	76
最低点	29	最低点	32
E大			
被験者数	92	被験者数	90
平均点	37.86	平均点	39.58
標準偏差	8.85	標準偏差	8.84
最高点	67	最高点	64
最低点	19	最低点	18
F大(5年)			
被験者数	26	被験者数	24

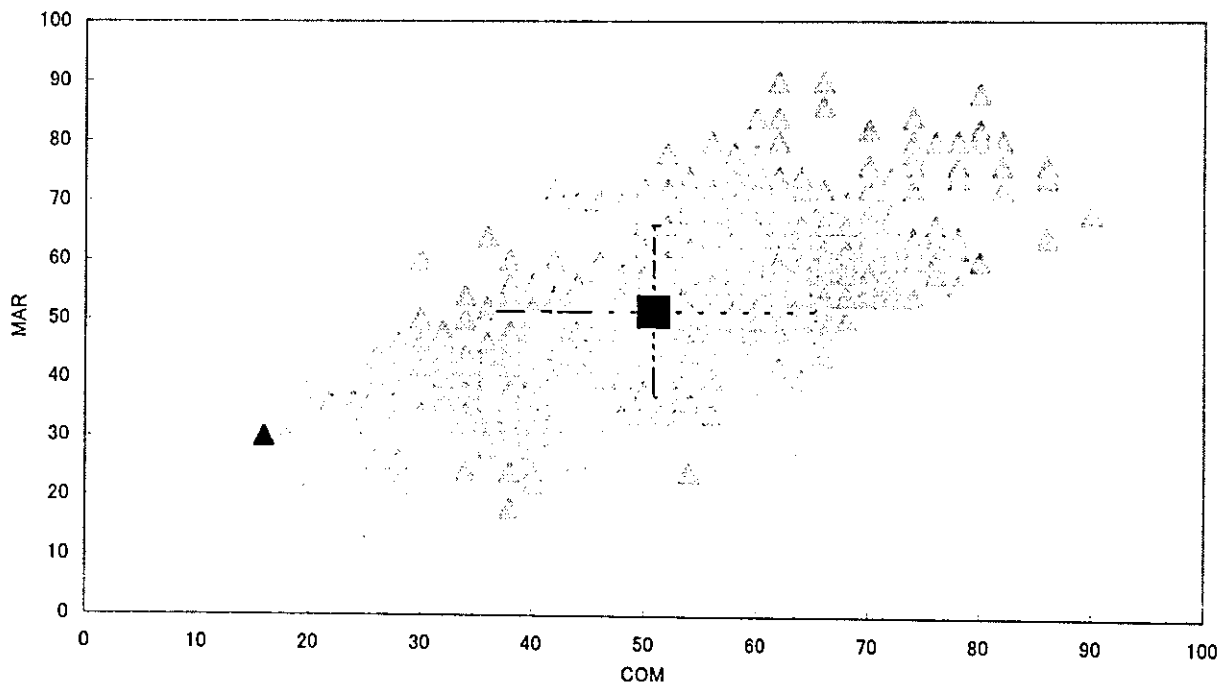


平均点	44.87	平均点	45.58
標準偏差	8.96	標準偏差	10.37
最高点	63	最高点	70
最低点	25	最低点	30

F大(6年)

被験者数	59	被験者数	58
平均点	66.74	平均点	65.90
標準偏差	11.31	標準偏差	10.87
最高点	92	最高点	88
最低点	40	最低点	36

(図1) MCQにおけるCOMとMARの得点比  
R=0.70



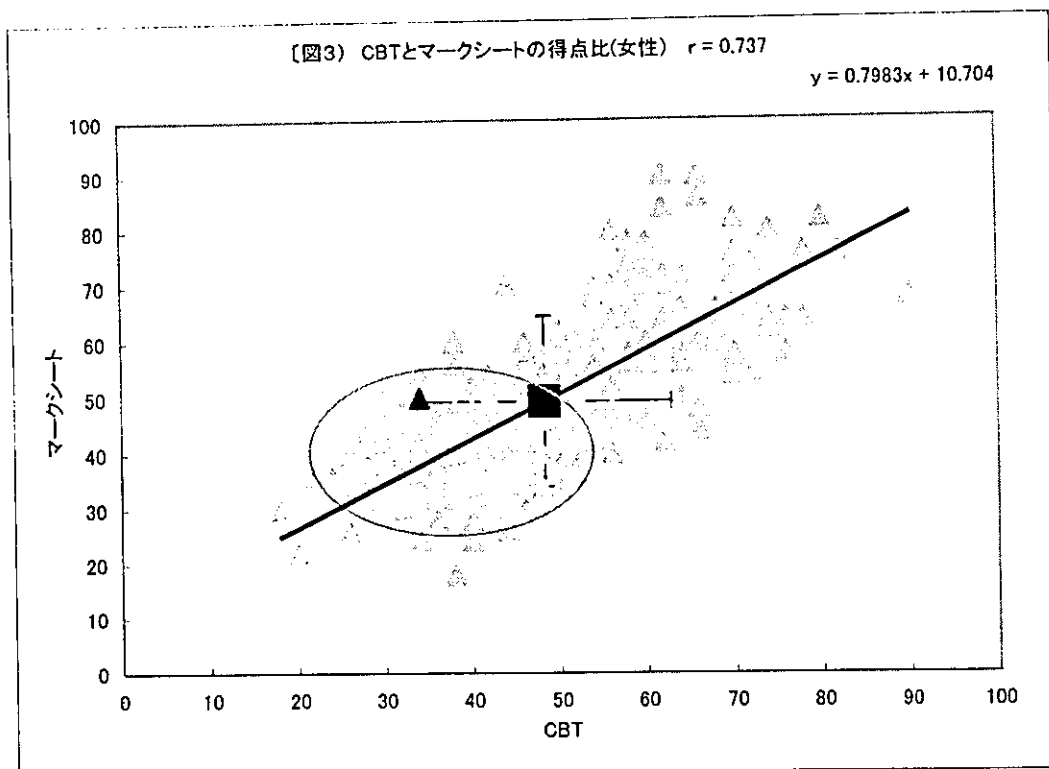
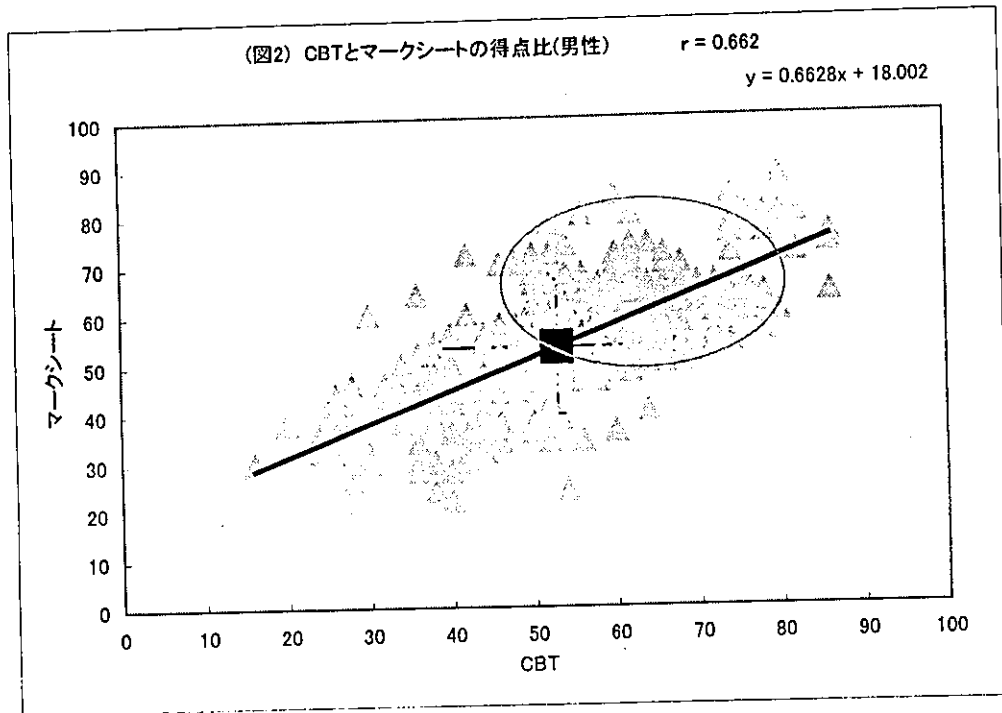
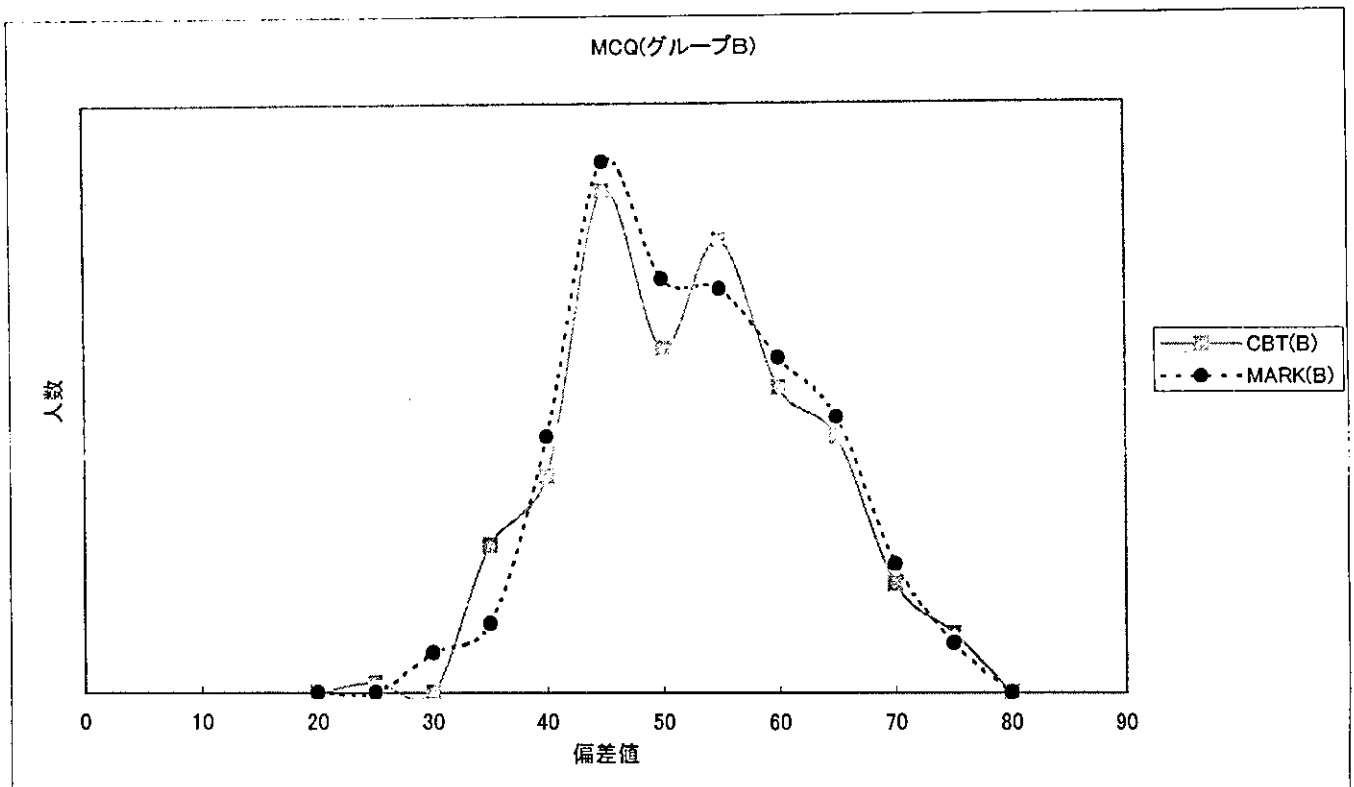
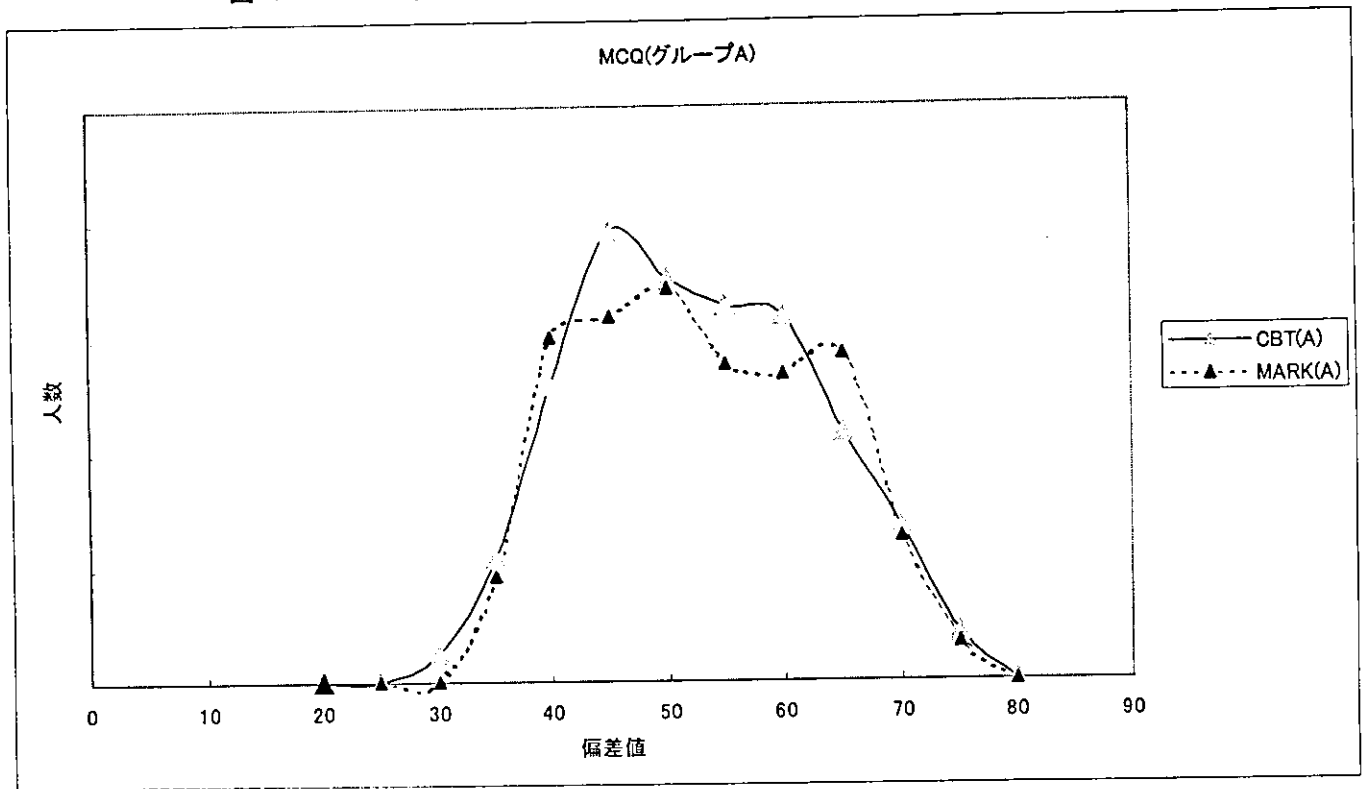
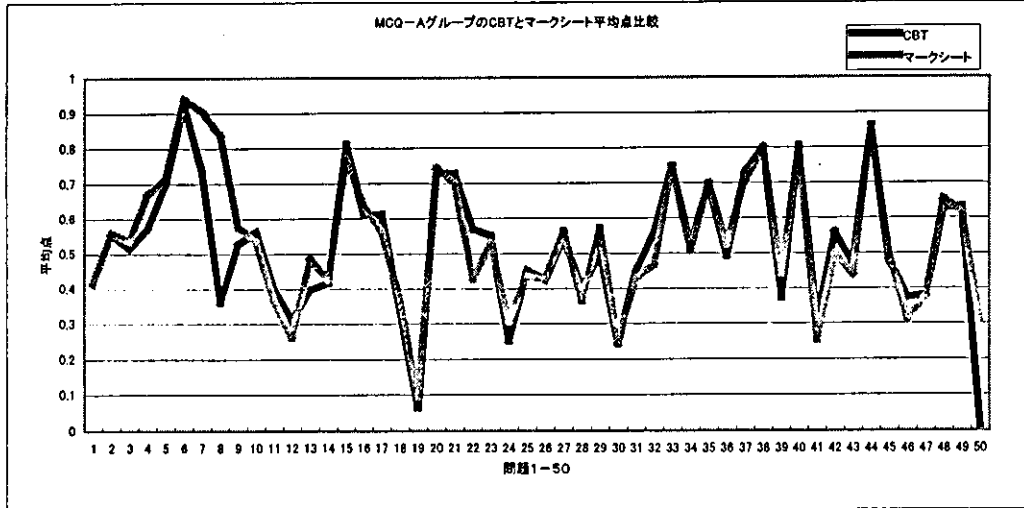
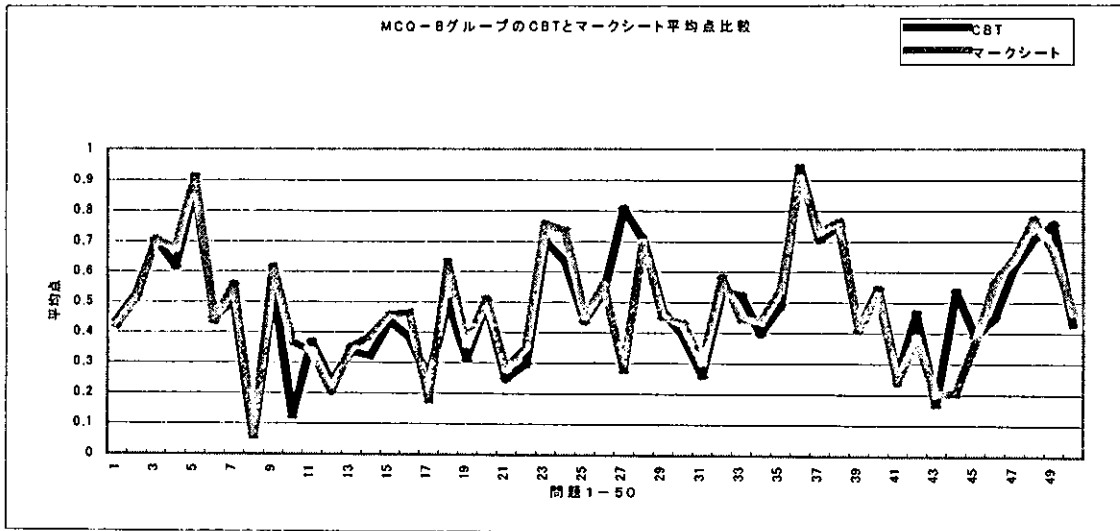


図4 COM (CBT) とMARにおける各問題の正答率

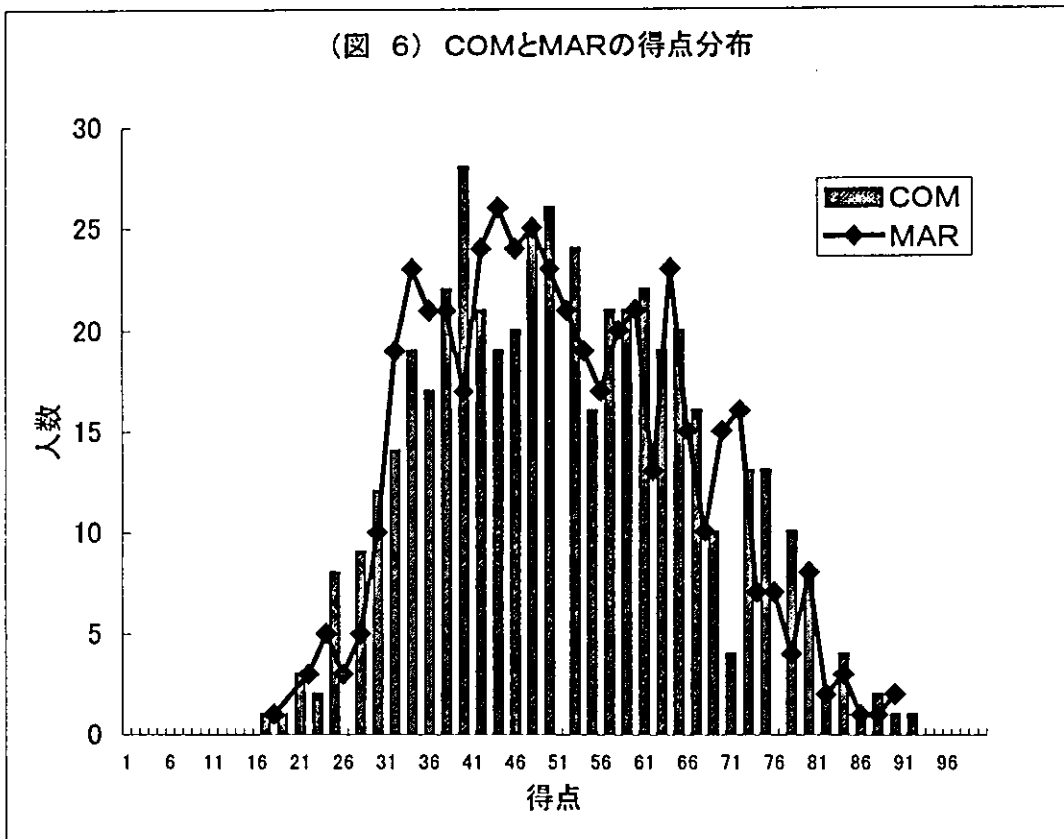




(図 5 A)



(図 5 B)



(表 2) PMP 問題の実施結果

C B T

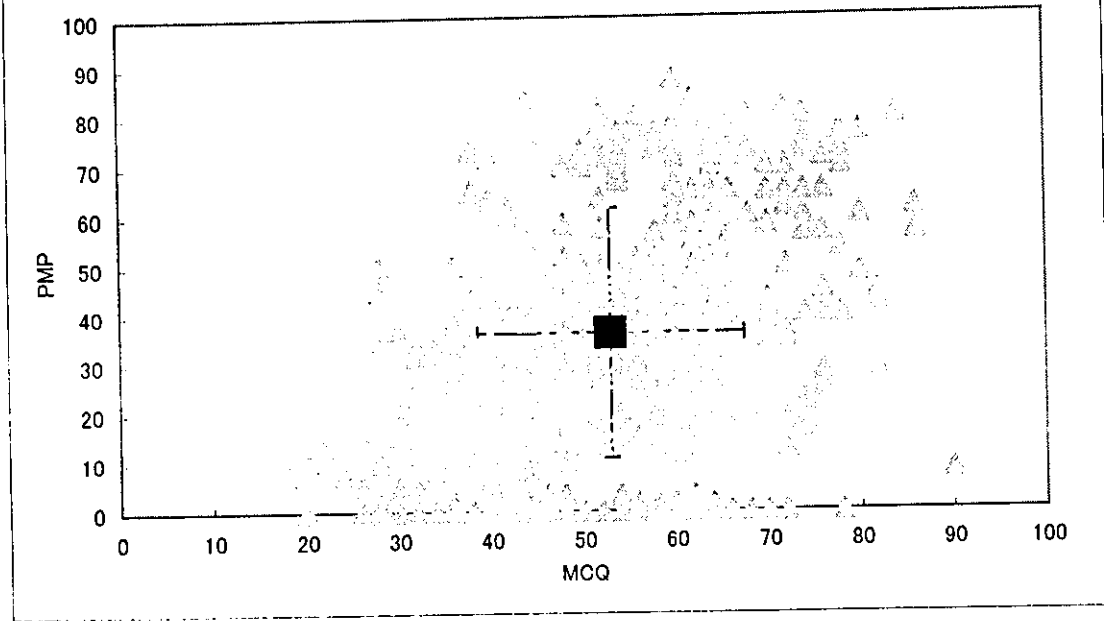
全体	アニサキス	喘息	肺梗塞	イレウス	過換気
被験者数 505	被験者数 244	被験者数 255	被験者数 240	被験者数 231	被験者数 485
平均点	平均点 61.27	平均点 33.82	平均点 29.40	平均点 29.87	平均点 31.8
標準偏差	標準偏差 45.43	標準偏差 27.49	標準偏差 20.74	標準偏差 31.41	標準偏差 37.5
最高点	最高点 100	最高点 95	最高点 85	最高点 100	最高点 100
最低点	最低点 0	最低点 0	最低点 0	最低点 0	最低点 0

MAR

全体	アニサキス	喘息	肺梗塞	イレウス
被験者数	被験者数 254	被験者数 244	被験者数 254	被験者数 244
平均点	平均点 60.39	平均点 51.39	平均点 74.96	平均点 54.10
標準偏差	標準偏差 19.76	標準偏差 22.14	標準偏差 22.42	標準偏差 20.88
最高点	最高点 100	最高点 100	最高点 100	最高点 100
最低点	最低点 0	最低点 0	最低点 0	最低点 0

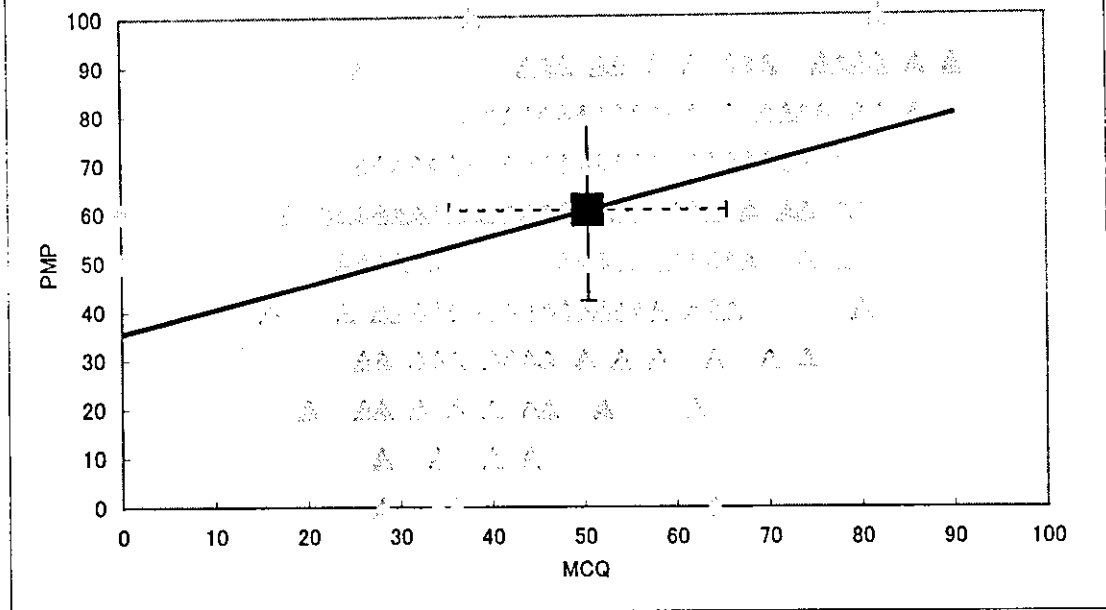
(図 7) コンピュータによるMCQ得点とPMP得点の分布

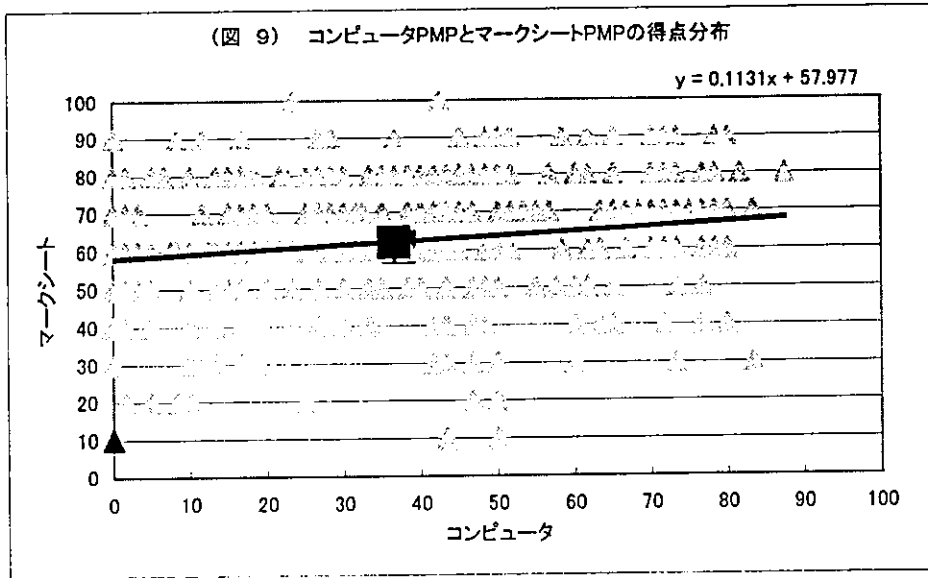
R=0.438845



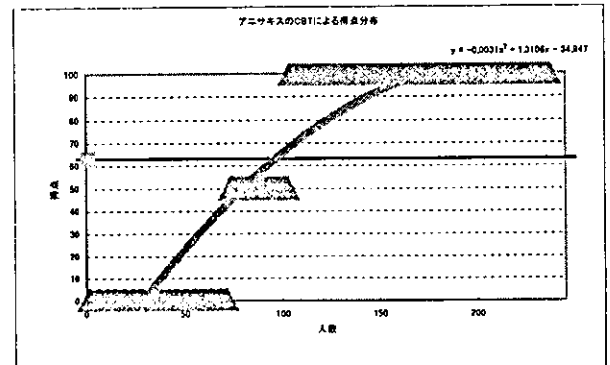
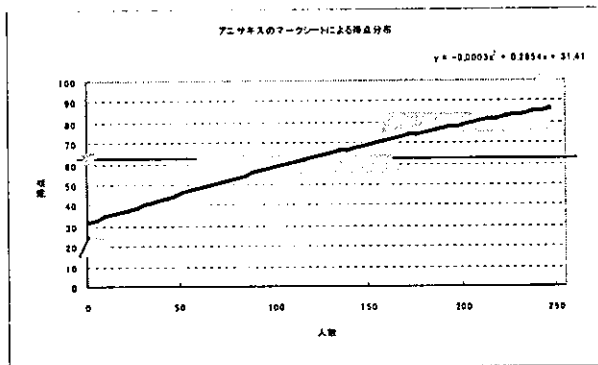
(図 8) コンピュータによるMCQ得点とマークシートPMP得点の分布

$y = 0.4949x + 35.49$

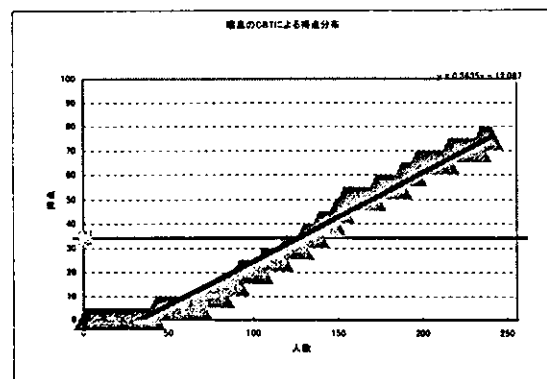
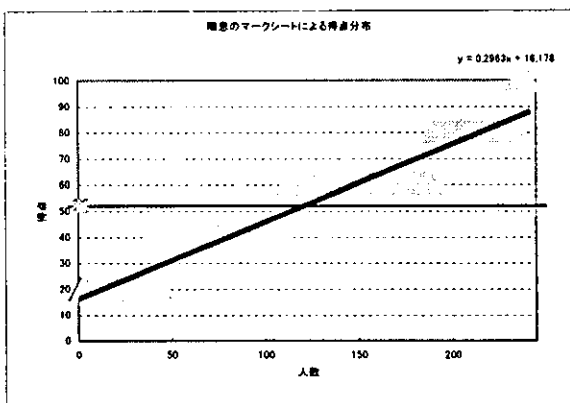




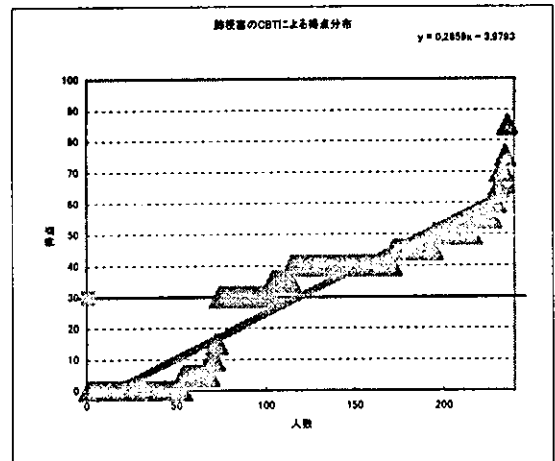
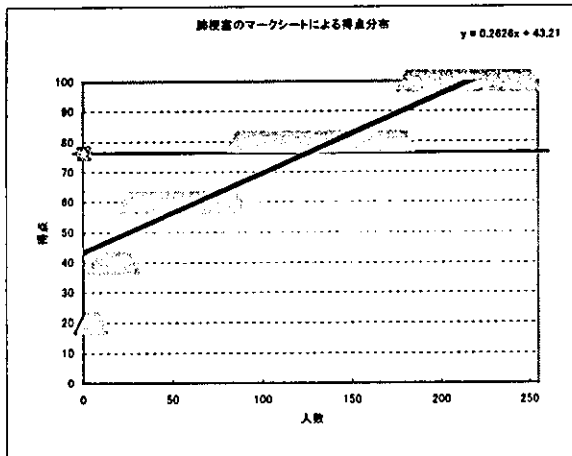
(図 10) アニサキスの問題の得点比較 (左がマークシート、右が PMP)



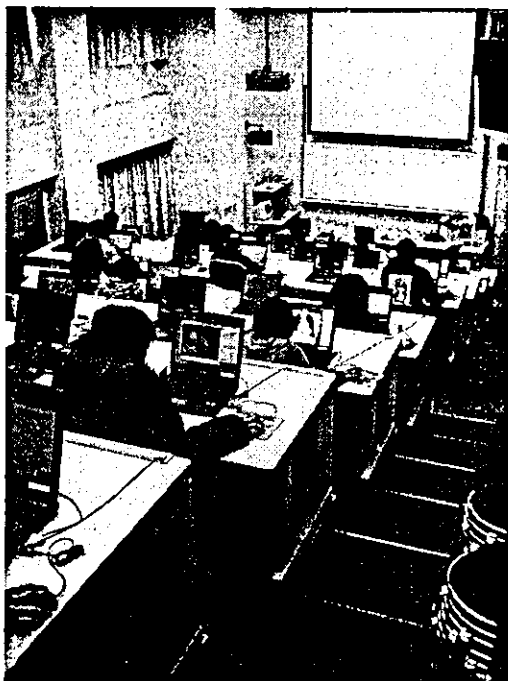
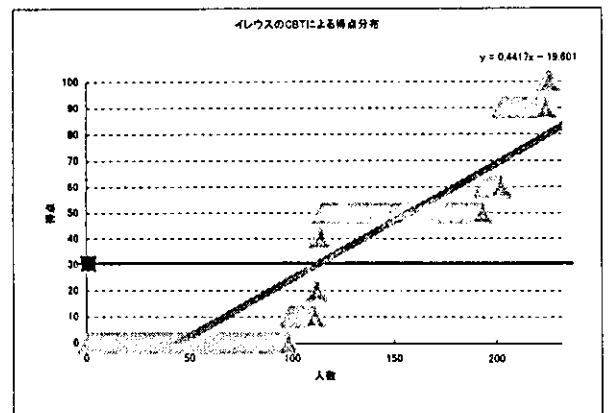
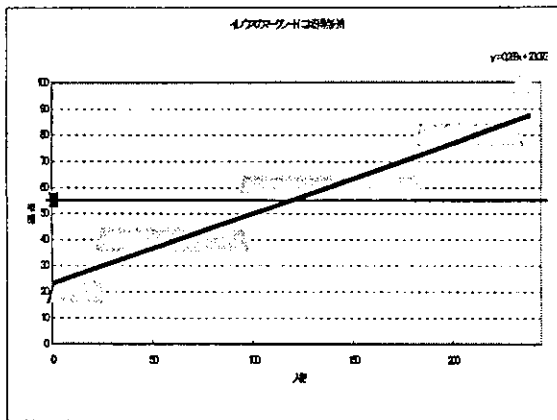
(図 11) 気管支喘息の問題の得点比較 (左がマークシート、右が PMP)



(図 1 2) 肺梗塞の問題の得点比較 (左がマークシート、右が PMP)



(図 1 3) イレウスの問題の得点比較 (左がマークシート、右が PMP)



(図 1 4) 階段教室の試験風景



分担研究報告書

試験問題のコンピュータ化と諸外国の実情視察

分担研究者 高林克日己 千葉大学医学部医療情報部助教授

研究要旨

研究班で医師国家試験の模擬問題をA、B 60題ずつ作成し、それぞれ筆記試験またはコンピュータ使用の試験として学生に実施し、その結果をまとめて比較し、コンピュータ化の実施可能性と利点、欠点について検討した。

A. 研究目的

医師国家試験にコンピュータを利用した試験を導入する妥当性を検討する。

B. 研究方法

①福島統分担研究者が作成した50問の試験問題2セット、計100問をコンピュータ化して、画像などとセットにし、ウインドウズマシンにインストール可能とした。同様に一昨年作成したツールを介してPMP問題を作成し、これらMCQ、PMPとアンケートをまとめてインストールできるように設計した。試験中には受験者の解答のほか、試験時間、とくにPMPでは各セクションに関する解答時間を収集し、これらをまとめて、ハードディスク上に一意の名前のファイル作成するようにした。受験者にはデータベースや画像は受験中に侵入できたり閲覧できないように設計した。実際のファイルの展開を示す。図1はMCQの問題の1例である。画像はクリックすることで大きくより鮮明な画像を表示できる。どの問題を終了したかは一覧でみることができまたどれでも番号を選ぶことでその問題に移ることができる。1選択、でも複数選択でも入力が可能である。また試験終了時にはすぐその場で評価点が示される。

PMP問題ははじめのオープニングシーンの情報〔図3〕以降は全て受験者が選択するもので、その結果は「カルテ」に貯まっていく〔図4〕。必要に応じて画像が展開することもある。最終的に

は診断と治療が行なわれ、これらの点数で評価が行なわれる。終了時には選択結果と各セクションの所要時間が表示される〔図5〕。これらを5問作成した。アンケートもこれと同様にコンピュータから行なった

②受験者のアンケートをとって情報を収集する。学生アンケートは全員の学生に試験終了時にコンピュータへの入力として行い、複数回答を許した。

C. 結果

①総括に付す

②1) MCQについて

表1、図6にPMP試験に対する学生の評価を示した。画像がみやすい、操作が楽という他に、紙よりよいとするものが20件あった。いっぽうでnegativeな意見としてメモが取りにくい、長文を読みにくい、難しかった、紙の試験の方がよい、目が疲れるという意見であった。

2) PMPについて

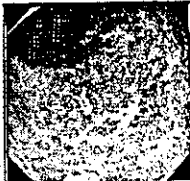
PMPについてのアンケートでは面白いとするものが賛成意見のほとんどを占めるのに対し、操作性の悪さ、難しさを訴えるものが多かった（表2、図7）。

D. 考察

①②総括に付す

Multi Choice Question

22歳の男性。1か月前から血性下痢を認め、昨日から腹部全体の痛みと38.5℃の発熱が出現し、鮮血の混入した下痢が1日20回以上となったため来院した。緊急下部内視鏡検査所見を別に示す。この疾患について正しいのはどれか。



設問1  設問2  設問3  設問4  設問5  
 設問6  設問7  設問8  設問9  設問10  
 設問11  設問12  設問13  設問14  設問15  
 設問16  設問17  設問18  設問19  設問20  
 設問21  設問22  設問23  設問24  設問25  
 設問26  設問27  設問28  設問29  設問30  
 設問31  設問32  設問33  設問34  設問35  
 設問36  設問37  設問38  設問39  設問40  
 設問41  設問42  設問43  設問44  設問45  
 設問46  設問47  設問48  設問49  設問50

1: 炎症は大腸壁全層に及ぶ。  
 2: 合併症として穿孔が多い。  
 3: 10年以内に悪化するものが多い。  
 4: サラソスルファピリジンが有効である。  
 5: Clostridium difficileの菌毒素が関与する。

1/50    C 1    C 2    C 3    C 4    C 5    答える    前へ    次へ    提出

図1 MCQの画面

結果

# 試験結果

0.0 %

了解

図2 試験終了時の表示

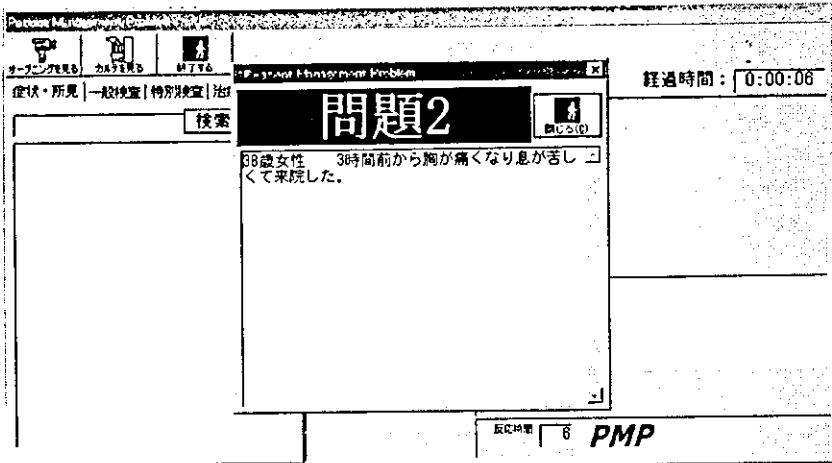


図3 PMP のオープニングシーン

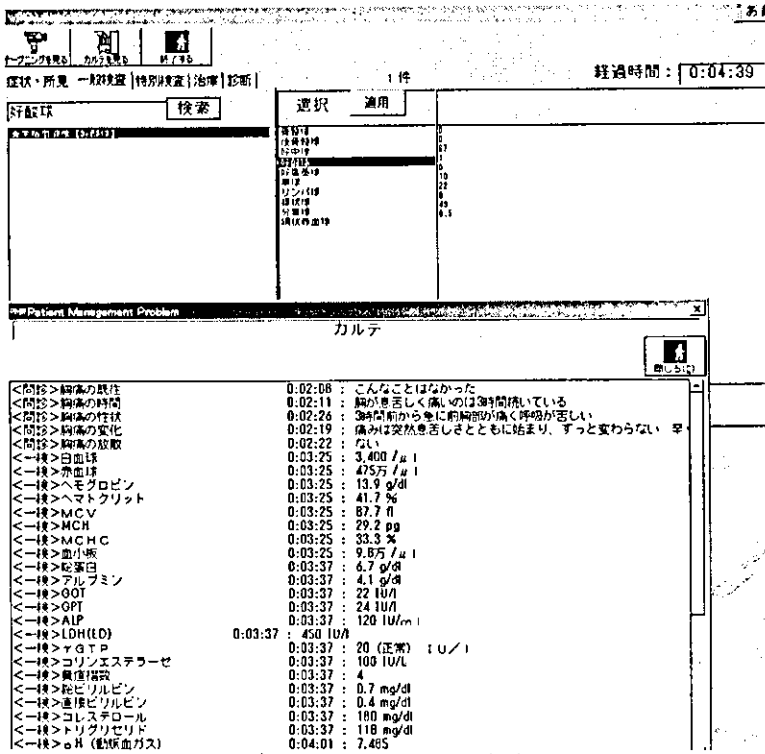


図4 PMP の項目選択の場面とカルテ内容

Patient Management Problem

症状・所見 | 一般検査 | 特別検査 | 治療 | 診断 | 1件 | 経過時間: 0:13:05

皮膚筋炎 | 検索 | 選択 | 適用

Patient Management Problem

履歴	問診	検査	治療	診断
	0:03:12		0:09:56	0:10:08

最大反応: 15 | 所要時間: 0:10:26

<問診> 胸痛の既往 0:02:08 : こんなことはなかった  
 <問診> 胸痛の時間 0:02:11 : 胸が息苦しく痛いのはず時間続いている  
 <問診> 胸痛の性状 0:02:26 : 3時間前から急に前胸部が痛く呼吸が苦しい  
 <問診> 胸痛の変化 0:02:19 : 痛みは突然息苦しさとともに始まり、ずっと変わ  
 らない、辛くて動かない  
 <問診> 胸痛の放散 0:02:22 : ない  
 <検査> 白血球 0:03:25 : 3,400 /μl  
 <検査> 赤血球 0:03:25 : 475万 /μl  
 <検査> ヘモグロビン 0:03:25 : 13.9 g/dl  
 <検査> ヘマトクリット 0:03:25 : 41.7 %  
 <検査> MCV 0:03:25 : 87.7 fl  
 <検査> MCH 0:03:25 : 29.2 pg  
 <検査> MCHC 0:03:25 : 33.3 %  
 <検査> 血小板 0:03:25 : 9.8万 /μl  
 <検査> 総蛋白 0:03:37 : 6.7 g/dl  
 <検査> アルブミン 0:03:37 : 4.1 g/dl  
 <検査> GOT 0:03:37 : 22 IU/l  
 <検査> GPT 0:03:37 : 24 IU/l  
 <検査> ALP 0:03:37 : 120 IU/ml  
 <検査> LDH(LD) 0:03:37 : 450 IU/l  
 <検査> γGTP 0:03:37 : 20 (GE常) IU/l  
 <検査> コリンエステラーゼ 0:03:37 : 100 IU/l  
 <検査> 糞便潜血 0:03:37 : 4  
 <検査> 総ビリルビン 0:03:37 : 0.7 mg/dl

図5 終了時 全ての設問に対する回答結果と時間が表示される