

XVII. Evaluarea testelor diagnostice: criterii de validitate

Anormalitatea

Medicii petrec o bună parte a timpului distingând normalul de anormal (Sunt aceste adenopatii patologice? Este consistența ficatului crescută? Este scăderea ponderală la acest individ determinată de stress sau are cancer?).

Când anomaliile sunt grosiere (adenopatii mari, ficat dur cu marginea ascuțită, pacient emaciat cu alte modificări biologice), este simplu de răspuns la aceste întrebări. Greutatea apare atunci când modificările sunt subtile, cum se întâmplă de obicei la debutul bolii – și uneori diagnosticul este cu atât mai important. Atunci fie decidem că modificările reflectă o boală și anunțăm pacientul, continuăm investigațiile (cu unele mai invazive și/sau mai scumpe) sau tratăm boala (operăm, administrăm tratamente care costă și au efecte adverse), fie decidem că modificările nu ies din sfera normalului și liniștim pacientul (dar ne asumăm riscul de a trece pe lângă o boală care totuși există, iar când diagnosticul va fi clar, poate va fi prea târziu).

Testele diagnostice sunt de două feluri, calitative și cantitative. Primele clasifică pacienții ca bolnavi sau sănătoși în funcție de prezența sau absența unei modificări clinice sau de laborator (de exemplu, prezența sau absența unei pneumonii la examenul radiologic). Cele cantitative au ca rezultat o variabilă numerică continuă, iar clasificarea pacientului ca bolnav sau nu se face pe baza unei valori de prag (numite criteriu de pozitivitate), în funcție de care testul este considerat pozitiv sau negativ (de exemplu TA, creatinina, transaminazele, VEMS etc.). Sunt mai multe metode de a hotărâ care este normalul:

1. **Metoda distribuției Gaussiene**, bazată pe presupunerea că valorile testului au o distribuție normală (Gaussiană). Conform acestei metode, se încadrează în normal cei aflați în zona mediei ± 2 deviații standard, adică 95% din populație, extremele fiind anormale. Pentru a obține aceste valori, se face un studiu pe o populație considerată “sănătoasă” (eventual separat pe grupe de vârstă, sex, rasă sau alte subgrupuri), se aplică testul diagnostic și se calculează media și deviația standard. Problema cu această metodă este aceea că rezultatele majorității testelor nu au o distribuție normală (de exemplu media creatininemiei este în jurul valorii 0,9, în timp ce extremele se pot duce și la 15 în sus, dar mult mai puțin în jos); de asemenea, conform acestei metode extremele de 2,5% ar reprezenta anormalul, deci toate bolile ar avea o prevalență fixă de 5%, ceea ce este absurd. Singurii pacienți normali ar fi cei necăutați (dacă aplicăm 1 test, probabilitatea să iasă un rezultat normal ar fi 95%; dacă aceluiași pacient îi aplicăm 2 teste, probabilitatea să iasă normale scade la 90%; iar dacă îi aplicăm 20 de teste independente – nu știu câte teste are screeningul anual obligatoriu introdus recent de Ministerul Sănătății – probabilitatea scade la $0.95^{20}=33\%$).
2. **Metoda procentuală**, care este aproape identică celei Gaussiene (se măsoară valorile testului într-o populație și se definește normalul tot ca valorile în care se încadrează 95% din populație, iar patologicul fie 5% din valorile superioare, fie 2,5% din extremele de ambele părți). Această metodă nu cere ca valorile testului să aibă o distribuție normală, dar are aceleași dezavantaje ca și metoda Gaussiană.
3. **Metoda “preferinței culturale”**, care descrie ca normal ceea ce societatea consideră astfel (de exemplu greutatea la femei). În acest caz se creează o confuzie cu privire la rolul medicinei.

4. **Metoda factorului de risc:** valorile devin patologice peste pragul la care este demonstrat că factorul respectiv devine factor de risc. Desigur, nu ne ajută cu nimic diagnosticarea unui factor de risc asupra căruia nu putem interveni.
5. **Metoda diagnostică sau a valorii predictive:** de această metodă ne vom ocupa în continuare, și pe aceasta o aplicăm de fapt în clinică zi de zi.
6. **Metoda terapeutică:** valorile devin patologice de la nivelul la care s-a demonstrat că tratându-le se obține mai mult beneficiu decât efecte adverse. Exemple pentru metoda aceasta și cea a factorului de risc pot fi TA și colesterolul, la care vedem că valorile de prag se modifică continuu. Pentru stabilirea valorilor de prag sunt necesare mari studii de cohortă (metoda factorului de risc) și terapeutice (metoda terapeutică). Metoda terapeutică este cea mai pragmatică, derivând din metoda factorului de risc.

Performanța măsurătorilor: validitatea și reproductibilitatea

Validitatea (sau **acuratețea**) este corectitudinea măsurătorii – ne arată cât de aproape este valoarea unei măsurători de valoarea reală (cât de aproape este o potasemie măsurată de potasemia reală).

Reproductibilitatea (sau **precizia**) arată gradul în care o serie de măsurători fluctuează în jurul unei valori centrale, valoare care poate fi mai mult sau mai puțin apropiată de valoarea reală (în funcție de **acuratețea** testului). De exemplu în determinarea unei potasemii (valoare reală de 4 mEq/l), metoda noastră de măsurare este validă dacă repetând măsurătoarea vom avea valori foarte apropiate de 4 mEq/l. Pentru ca metoda să fie precisă, măsurătorile trebuie să dea valori foarte apropiate între ele, chiar dacă diferite de cea reală (de exemplu dacă toate măsurătorile dau valoarea 7,49 mEq, înseamnă că metoda este foarte precisă, dar deloc validă).

Acuratețea se poate evalua determinând diferența dintre media rezultatelor măsurării și valoarea reală, în timp ce precizia se evaluează observând distribuția frecvențelor măsurătorilor și calculând deviația standard a acestora. O metodă de măsurare este bună atunci când este în același timp validă și precisă.

Atunci când rezultatele testelor sunt binare (pozitiv/negativ), reproductibilitatea interobservator se calculează raportând numărul de rezultate convergente (în care investigatorii au dat același verdict) la numărul de rezultate divergente (în care părerile investigatorilor au fost diferite) - aceasta este reproductibilitatea calculată prin metoda simplă (Tabelul XVII.1 - doi clinicieni au consultat 100 de pacienți cu dispnee: ambii sunt de acord cu prezența zgomotului 3 la 5 pacienți și cu absența lui la 75 de pacienți – așadar, “acordul simplu” este de $(5+75)/100 = 0,80$, sau acordul a fost prezent în proporție de 80%).

Tabelul XVII.1 Calculul reproductibilității interobservator – metoda simplă. Acordul dintre doi clinicieni asupra prezenței galopului protodiastolic la pacienți cu dispnee.

		Clinician II (galop)	
		DA	NU
Clinician I (galop)	DA	5	5
	NU	15	75

Acest calcul are avantajul de a fi simplu de efectuat, dar nu ține cont de faptul că cei doi observatori pot ajunge la același rezultat, măcar într-un număr de cazuri, din întâmplare (să ne închipuim că medicii, în loc să asculte zgomotele cardiace, diagnostichează prezența zgomotului 3 dând cu banul: capul = zgomot 3 prezent, pajura = zgomot 3 absent; este clar că, de câte ori cei doi vor ajunge la același rezultat prin această “metodă diagnostică”, acest lucru va fi rezultatul întâmplării).

Pentru remedierea acestui neajuns, deci pentru a vedea care este acordul dintre doi investigatori în realitate, scăzând contribuția întâmplării, se folosește **coeficientul de concordanță κ** , care ia valori de la 0 la 1 (pentru valorile de la 0 la 0,2 concordanța este ușoară, de la 0,2 la 0,4 este acceptabilă, de la 0,4 la 0,6 moderată, de la 0,6 la 0,8 substanțială, iar de la 0,8 la 1 aproape perfectă¹). Pentru exemplul din Tabelul XVII.1, coeficientul de concordanță κ este 0,23 (mult mai mic decât acordul simplu de 80%, considerându-se că acordul datorat întâmplării a fost în acest caz 74% - pentru calculul κ , vezi^{1, 2}).

În Tabelele XVII.2 și 3 putem vedea câteva exemple de acord interobservator (pentru o colecție mult mai mare, vezi² și mai ales¹). Se observă că testele diagnostice pe care le folosim (de la semne clinice la teste paraclinice) sunt departe de a fi perfecte în privința reproductibilității. Chiar și în privința analizelor de laborator, în care clinicianul nu are decât de interpretat un număr, dezacordul dintre investigatori este încă posibil (de exemplu, într-un studiu în care trei endocrinologi au văzut aceleași rezultate privind funcția tiroidiană și alte date clinice privind 55 de pacienți consecutivi la care se suspecta boală tiroidiană, ei au fost în dezacord cu privire la diagnosticul final în 40% din cazuri³. Mai mult, nici măcar analiza computerizată a rezultatelor nu are o reproductibilitate mai bună: într-un studiu asupra unor perechi de electrocardiograme făcute la interval de un minut la 92 de pacienți, interpretarea computerului a fost semnificativ diferită în 40% din cazuri, chiar dacă traseele nu difereau⁴.

Tabelul XVII.2. Acordul interobservator – câteva exemple în cazul semnelor fizice (din¹).

Semn	Coeficient κ
Piele	
Paloare (pacientul pare anemic) ^{5, 6}	0,23-0,48
Paloare conjunctivală ⁷	0,54-0,75
Cianoză ^{5, 8}	0,36-0,70
Icter ⁹	0,65
Semne vitale	
Hipotensiune (TA sistolică < 90mmHg)	0,90
Febră (evaluată prin palparea pielii) ⁵	0,09-0,23
Tahipnee ⁸	0,25
Retinopatie diabetică^{10, 11}	
Exudate	0,56-0,67
Hemoragii intraretiniene	0,89
Neovascularizație	0,1-0,48

Stadializare	0,65
Matitate la percuția pulmonară ^{8,12, 13}	0,16-0,52
Raluri alveolare ^{8, 12, 14}	0,21-0,63
Raluri sibilante ^{8, 12, 14, 15}	0,43-0,93

Tabelul XVII.3. Acordul interobservator – câteva exemple în cazul testelor diagnostice (din¹).

Semn (la testul diagnostic)	Coefficient κ
Radiografie toracică	
Cardiomegalie ¹⁶	0,48
Redistribuirea circulației pulmonare ¹⁶	0,50
Fibroza pulmonară (pe o scală cu 4 grade) ¹⁷	0,45
Venografie cu substanță de contrast	
Tromboză venoasă profundă ¹⁸	0,53
Angiografie cu subtracție digitală	
Stenoză de arteră renală ¹⁹	0,65
Coronarografie	
Clasificarea leziunilor arterelor coronare ²⁰	0,33
Tomografia computerizată cerebrală²¹	
Normală sau anormală, la pacienți cu accident vascular cerebral	0,60
Leziune pe partea dreaptă sau stângă, la pacienți cu accident vascular cerebral	0,65
Efect de masă, prezent sau absent	0,52
Tomografia computerizată toracică	
Stadializarea cancerului pulmonar ²²	0,40-0,60
Rezonanță magnetică nucleară cerebrală	
Compatibilă cu scleroză multiplă ²³	0,57-0,87
Rezonanță magnetică nucleară a coloanei vertebrale	
Modificări ale discului intervertebral sau normal ²⁴	0,59
Ecografie	
Tromboză venoasă profundă, prezentă sau absentă ²⁵	0,69
Nodul tiroidian, prezent sau absent ^{26, 27}	0,57-0,66
Examenul histopatologic al biopsiei hepatice²⁸	
Colestază	0,40
Boală alcoolică a ficatului	0,49
Ciroză	0,59

Prin efectuarea unui test, clinicianul se străduiește să clasifice starea reală dar necunoscută a unui subiect de observație cu ajutorul unui instrument imperfect.

Sursele de nesiguranță/variabilitate sunt:

1. Instrumentul de măsură: imprecizia analitică (aceiași test aplicat aceluiași pacient să dea același rezultat).
2. Variabilitatea subiectului: intraindividuală (regresia către medie – fluctuații fiziologice)/ interindividuală (cu cât această variabilitate este mai mare, avem nevoie de un eșantion mai mare pentru studiul nostru).
3. Variabilitatea interpretării: interindividuală (doi radiologi sau ecografiști văd lucruri diferite la același pacient) / intraindividuală (aceiași radiolog/anatomopatolog vede lucruri diferite citind aceeași radiografie/lamă în momente diferite; studiile au arătat că aceste diferențe sunt incredibil de mari). Această variabilitate este cuantificată prin mărimea coeficientului de concordanță κ .
4. Validitatea intrinsecă a testului (sensibilitatea, specificitatea, raportul de probabilitate=*likelihood ratio*).
5. Prevalența bolii = probabilitatea pretest, care influențează valorile predictive ale testului conform teoremei lui Bayes.

Ca la orice tip de studiu, și la cele terapeutice ne interesează două lucruri: validitatea (calitatea metodologică) și rezultatele. Este necesar ca în primul rând studiul să fie corect efectuat, pentru a ne putea baza pe rezultate.

Evaluarea validității studiului se face prin verificarea următoarelor criterii:

1. A fost testul comparat cu un *gold standard* adevărat? A fost comparația “oarbă”?

Gold standard-ul este testul etalon, cu care comparăm orice test nou și în funcție de care îl evaluăm pe acesta. Se presupune că *gold standard*-ul este testul perfect, care identifică toți indivizii care au boala și nu dă rezultate fals pozitive sau negative – de obicei este examenul histopatologic (în coronaropatii, este coronarografia). Există boli pentru care nu avem *gold standard* – de exemplu bolile pentru care avem criterii de diagnostic. În aceste cazuri trebuie căutat un alt tip de *gold standard*, cum ar fi evoluția clinică (pentru artrita reumatoidă, de exemplu, pentru care nu avem un test perfect). Deloc surprinzător, de cele mai multe ori tocmai pentru astfel de boli fără *gold standard* se caută un test diagnostic bun. Acolo unde avem *gold standard* evaluăm teste noi fie pentru că sunt mai ieftine, fie mai puțin invazive, fie mai ușor de efectuat. Atunci când testul etalon este imperfect, și evaluăm un test nou în funcție de acesta, dacă testul nou este mai bun, prin natura studiilor diagnostice acest fapt nu va putea fi evidențiat, iar testul nou întotdeauna apare ca fiind mai slab decât *gold standard*-ul.

Testul de evaluat trebuie să fie independent de *gold standard* – de exemplu, nu putem evalua anticorpii antiADN dublu catenar ca test diagnostic în lupus, având ca *gold standard* criteriile ACR pentru lupus, deoarece printe criterii se află și acești autoanticorpi. Iar la evaluarea criteriilor de diagnostic pentru arterita cu celule gigante²⁹, *gold standard*-ul folosit a fost... părerea experților. Oare pe ce s-o fi bazat părerea acestora, dacă nu (măcar parțial) tot pe simptomele și testele cuprinse în aceste criterii! Așadar, atenție la ce *gold standard* a fost folosit, și cât de valid este acesta, pentru a aprecia validitatea studiului!

De asemenea, comparația dintre *gold standard* și testul de evaluat trebuie să fie oarbă – adică cei care au efectuat și interpretat testul de evaluat nu trebuie să știe rezultatul *gold*

standard-ului și invers. După ce medicii află de existența unui nodul pulmonar apărut la CT, îl vor vedea și pe radiografie, iar după ecocardiografia care a arătat o regurgitare aortică, normal că vor auzi și suflul! Astfel se introduce o eroare sistematică în favoarea testului de evaluat. Bineînțeles că, la fel ca la studiile terapeutice, cu cât testul este mai subiectiv (examen fizic, anamneză, radiologie, ecografie, chiar examen histologic – am vazut cât de mare este variabilitatea inter/intraobservator), cu atât mai mult este nevoie de precauții care să asigure orbirea. Pentru testele biochimice – care nu pot fi influențate de cel care le interpretează – lipsa orbirii nu prea poate introduce erori sistematice în favoarea acurateții testului³⁰ - deși nu este imposibilă interpretarea subiectivă a rezultatelor biochimice³!

2. A fost testul evaluat la pacienții potriviți?

Pacienții pe care efectuăm studiul diagnostic trebuie să fie asemănători acelorora la care l-am folosi în practică. Un test diagnostic este util atunci când face diferența între afecțiuni asemănătoare, între care fără el nu am prea putea face diagnosticul diferențial. Oricine își dă seama că un pacient în vârstă, cu istoric de infarct anterior, ortopnee, subcrepitante, edeme, jugulare turgescențe, cardiomegalie și galop protodiastolic are insuficiență cardiacă, iar o tânără cu crize de wheezing are altă cauză a dispneei, nu este nevoie pentru asta să dozăm peptidul natriuretic. La fel, un test pentru artrita reumatoidă ne este necesar la pacienți care au poliartrită de scurt timp, și ne întrebăm dacă este vorba despre o artrită reumatoidă sau o altă boală inflamatorie articulară. Dacă, pentru a evalua un test diagnostic se folosesc, ca bolnavi, pacienți cu aspect clar al bolii respective, iar ca indivizi neavând boala oameni sănătoși, testul va apărea cu o putere discriminativă mult mai mare decât în realitate. Așadar, trebuie ca testul să fie evaluat la pacienți consecutivi la care boala respectivă este suspectată.

3. A fost efectuat *gold standard*-ul la toți pacienții, indiferent de rezultatul testului evaluat?

După cum am mai spus, presupunând că testul etalon este perfect, avem totuși nevoie să dezvoltăm alte teste, mai ieftine, mai puțin laborioase sau mai puțin invazive decât acesta. Nu putem face coronarografie la toți pacienții pentru că este o investigație invazivă și disponibilă în puține centre (la noi în țară), care necesită personal specializat, de asemenea preferăm să avem un alt test, imagistic prin care să diferențiem formațiunile intrahepatice decât examenul histopatologic al fragmentului obținut prin biopsie țintită. Din acest motiv căutăm să dezvoltăm alte teste, mai simple, cum ar fi computerul tomograf pentru bolile coronariene, leziunile colonice sau tromboembolismul pulmonar, RMN pentru afecțiunile coledocului - în loc de coronarografie, colonoscopie, scintigrafie/arteriografie, colangiografie retrogradă endoscopică etc. - testele etalon. În studiile care consacreză aceste noi investigații, pentru a putea face comparația testului nou cu *gold standard*-ul, trebuie ca tuturor pacienților să li se aplice concomitent și testul de evaluat, și *gold standard*-ul și nu, de exemplu, dacă pacientului nu-i iese nimic la colangio-RMN, să nu-i mai facem colangiografia retrogradă endoscopică, sau dacă tomografia cu emisie de pozitroni nu arată că leziunea este malignă, să nu-i mai facem biopsie din tumoră prin puncție sau prin operație, ori dacă anticorpii anti-peptid citrulinat (antiCCP) sunt negativi, să ne spunem că oricum nu părea o artrită reumatoidă și să nu mai urmărim pacientul un an (evoluția clinică fiind, în acest caz, *gold standard*-ul), deci să considerăm că pacientul nu are boala, fără să-i mai aplicăm *gold standard*-ul. Ori în toate aceste cazuri am omite cazurile fals negative și testul de evaluat va părea, din nou, mai bun, cu o sensibilitate mult mai mare decât în realitate. Invers, dacă pacienții la care

testul de evaluat iese pozitiv îi considerăm ca având boala, renunțând să le mai facem testul *gold standard*, oțitem fals negativii, și atunci testul va părea că are o specificitate mult mai mare decât în realitate.

Când pacienții au un test diagnostic negativ, investigatorii sunt tentați să nu mai aplice *gold standard*-ul, iar când acesta este invaziv sau riscant (de exemplu angiografie) poate că nici nu merită să fie făcut la pacienții la care testul de evaluat a fost negativ. Pentru a trece peste această deficiență, investigatorii pot utiliza în aceste cazuri un nou standard de referință pentru a demonstra că pacienții într-adevăr nu au avut boala – de exemplu evoluția clinică fără probleme în absența tratamentului. O dovadă convingătoare că un pacient cu suspiciune clinică de tromboză venoasă profundă nu a avut-o include lipsa oricărei complicații pe durata unei urmăriri îndelungate fără tratament anticoagulant³¹ - ceea ce, de fapt, nu constituie de loc o dovadă că pacientul nu a avut boala. În schimb, este o evaluare mult mai pragmatică a valorii unui test diagnostic – nici nu mă interesează să descopăr o tromboză care nu face rău, pentru că nu este nevoie să o tratez – așa s-a procedat și în celebrul studiu PIOPED³².

De obicei aplicăm testele pe bolnavi, și atunci, în lipsa unui studiu diagnostic bun, ne biziim pe aceste rezultate deformate de lipsa cunoștințelor privind rezultatele testului la pacienții care nu au boala. Rezonanță magnetică nucleară (RMN), de exemplu, se face pacienților cu dureri la nivelul coloanei lombare, și se descoperă tot felul de modificări, care au determinat explozia intervențiilor chirurgicale la acest nivel la începutul anilor 90. Întrebarea legitimă este: cum arată coloana la indivizi fără dureri, și pentru a răspunde s-a efectuat RMN la 98 de voluntari care nu se plâneau de nimic³³. Radiologii care au citit rezultatele nu știau asta (erau “orbi”), și au descoperit la fel de des protruzii discale ca și la pacienții care făceau RMN pentru dureri (și care fuseseră, probabil, operați de mult...).

4. S-a demonstrat că testul este reproductibil inter/intraobservator?

Dacă un medic efectuează un test de două ori asupra unui subiect a cărui condiție nu s-a schimbat, într-o anumită proporție de cazuri el va obține rezultate diferite – acest lucru este valabil pentru toate testele și toți investigatorii, numai că gradul de concordanță variază – una este să fie de 99%, alta 50%. Cu atât mai mult variază rezultatele când testul este efectuat de investigatori diferiți. Așadar, înainte de a începe evaluarea unui test, trebuie să vedem dacă reproductibilitatea lui este acceptabilă (nu este mai puțin adevărat că același lucru se întâmplă și cu *gold standard*-ul – am mai spus, există variabilitate intra- și interobservator și la citirea aceleiași lame histologice; în aceste cazuri se poate apela la citirea testului de către doi sau mai mulți investigatori, “orbi” unul față de celălalt).

5. Au fost furnizate intervalele de încredere pentru sensibilitate, specificitate și ceilalți parametri ai testului?

6. Este furnizat raportul de probabilitate (Likelihood ratio) al testului, sau datele din care acesta poate fi calculat?

Bibliografie

¹. McGee S. Evidence-based physical diagnosis. WB Saunders, Philadelphia, 2001. p.33-50.

2. Sackett DL, Haynes B, Guyatt G, Tugwell P. Clinical epidemiology. A basic science for clinical medicine. 2nd edition, Little, Brown. Toronto. 1991. p. 25-35.
3. Jarlov AE et al. Observer variation in the clinical and laboratory evaluation of patients with thyroid dysfunction and goiter. *Thyroid* 1998; 8:393-398.
4. Spodick DH, Bishop RL. Computer treason: intraobserver variability of an electrocardiographic computer system. *Am J Cardiol* 1997; 80:102-103.
5. Gjorup T et al. Global assessment of patients – a bedside study. II. Interobserver variation and frequency of clinical findings. *J Intern Med* 1990; 228:147-150.
6. Gjorup T et al. A critical evaluation of the clinical diagnosis of anemia. *Am J Epidemiol* 1986; 124:657-665.
7. Sheth TN et al. The relation of conjunctival pallor to the presence or absence of anemia. *J Gen Intern Med* 1997; 12:102-106.
8. Spiteri MA, Cook DG, Clarke SW. Reliability of eliciting physical signs in examination of the chest. *Lancet* 1988; 2:873-875.
9. Espinoza P et al. Interobserver agreement in the physical diagnosis of alcoholic liver disease. *Dig Dis Sci* 1987; 32:244-247.
10. Milton RC, Ganley JP, Lynk RH. Variability in grading diabetic retinopathy from stereo fundus photographs: comparison of physician and lay readers. *Br J Ophthalmol* 1977; 61:192-201.
11. Early Treatment Diabetic Retinopathy Study Research Group. Grading diabetic retinopathy from stereoscopic color fundus photographs: an extension of the modified Airlie House classification. ETDRS report number 10. *Ophthalmology* 1991; 98:786-806.
12. Mulrow CD et al. Observer variation in the pulmonary examination. *J Gen Intern Med* 1986; 94:188-196.
13. Gjorup T, Bugge PM, Jensen AM. Interobserver variation in assessment of respiratory signs: physicians' guesses as to interobserver variation. *Acta Med Scand* 1984; 216:61-66.
14. Holleman DR, Simmel DL, Goldberg JS. Diagnosis of obstructive airways disease from the clinical examination. *J Gen Intern Med* 1993; 8:63-68.
15. Badgett DG et al. Can moderate chronic obstructive pulmonary disease be diagnosed by historical and physical findings alone? *Am J Med* 1993; 94:188-196.
16. Butman SM et al. Bedside cardiovascular examination in patients with severe chronic heart failure: Importance of rest or inducible jugular venous distension. *J Am Coll Cardiol* 1993; 22:968-974.
17. Baughman RP et al. Crackles in interstitial lung disease: comparison of sarcoidosis and fibrosing alveolitis. *Chest* 1991; 100:96-101.
18. Illescas FF et al. Interobserver variability in the interpretation of contrast venography, technetium-99m red blood cell venography and impedance plethysmography for deep venous thrombosis. *J Can Assoc Radiol* 1990; 41:264-269.

-
- ¹⁹. DeVries ar et al. Interobserver variability in assessing renal artery stenosis by digital subtraction angiography. *Diagn Imag Clin Med* 1984; 53:277-281.
- ²⁰. Herman JPR et al. Inter- and intra-observer variability in the qualitative categorization of coronary angiograms. *Int J Card Imag* 1996; 12:21-30.
- ²¹. Shinar D et al. Interobserver reliability in the interpretation of computed tomographic scans of stroke patients. *Arch Neurol* 1987; 44:149-155.
- ²². Webb WR et al. Interobserver variability in CT and MR staging of lung cancer. *J Comput Assist Tomogr* 1993; 17:841-846.
- ²³. Barkhof F et al. Interobserver agreement for diagnostic MRI criteria in suspected multiple sclerosis. *Neuroradiology* 1999; 41:347-350.
- ²⁴. Jensen MC et al. Magnetic resonance imaging of the lumbar spine in eople without back pain. *N Engl J Med* 1994; 331:69-73.
- ²⁵. Atri M et al. Accuracy of sonography in the evaluation of calf deep vein thrombosis in both postoperative surveillance and symptomatic patients. *Aqm J Radiol* 1996; 166:1361-1367.
- ²⁶. Jarlov AE et al. Observer variation in ultrasound assessment of the thyroid gland. *Br J Radiol* 1993; 66:625-627.
- ²⁷. Schneider AB et al. Thyroid nodules in the follow-up of irradiated individuals: Comparison of thyroid ultrasound with scanning and palpation. *J Clin Endocrinol Metab* 1997; 82:4020-4027.
- ²⁸. Theodossi A et al. Observer variation in assessment of liver biopsies including analysis by kappa statistics. *Gastroenterology* 1980; 79:232-241.
- ²⁹. Hunder GG, Bloch DA, Michel BA et al. The American College of Rheumatology 1990 criteria for the classification of giant cell arteritis. *Arthritis Rheum.* 1990; 33:1122-8.
- ³⁰. Meissner K, Distel H, Mitzdorf U. Evidence for placebo effects on physical but not on biochemical outcome parameters: a review of clinical trials. *BMC Med.* 2007; 5:3
- ³¹. Strauss SE, Richardson S, Glasziou P, Haynes B. Evidence-based medicine. How to practice and teach EBM. 3rd edition, Elsevier, London, 2005. p.73.
- ³². The PIOPED investigators. Value of the ventilation/perfusion scan in acute pulmonary embolism: results of the prospective investigation of pulmonary embolism diagnosis (PIOPED). *JAMA* 1990; 263: 2753-59.
- ³³. Jensen MC, Brant-Zawadzki MN, obuchowski N, Modic MT, Malkasian D, Ross JS. Magnetic resonance imaging of the lumbar spine in people withot back pain. *N Engl J Med* 1994; 331: 69-73.