

1. ÖLÇEKLERDE GEÇERLİK VE GÜVENİRLİK

Dr.Demirali Y.Ergin

Her ölçme aracında vazgeçilemez iki özelliğin olması aranır: Bunlar geçerlik ve güvenilirliktir. Bir ölçme aracı, her şeyden önce ölçülecek özelliği tam ve doğru olarak ölçmeli, başka bir özellikle karıştırmadan ölçmelidir: ölçüğün bu niteliğine geçerlik denir Geçerlik, testin kullanılış amacına uygun hizmet etme derecesini belirler. Bir ölçme aracı ölçtüğü özelliği tutarlı olarak ölçebilmelidir. Aynı şartlar altında tekrarlandığında aynı sonuçları verebilmelidir. Bu ise güvenilirlik olarak adlandırılır.

Öncelikle, geçerlik ve güvenilirlik çalışmalarının hangi ölçme araçları için söz konusu olduğunu değerlendirmek gerekir. Her- ölçme aracı geçerlik ve güvenilirlik özelliğinin her ikisini birden taşımalıdır. Ama ölçme araçlarının ne ölçüde geçerli ve güvenilir' olduğunu belirlemek istatistik biliminden yararlanılarak hile olsa her zaman kolay ve hatta mümkün olmayabilir. Doğa bilimlerinde olduğu gibi olgusal verilerin elde edildiği ölçme araçları için geçerlik ve güvenilirlik hesaplanabilir ve hesaplanmalıdır. Özellikle psikolojik temelli bilimlerde olduğu gibi yargısal verilerin elde edildiği çileme araçlarında ise farklı durumlar ortaya çıkmaktadır', İster 'yapırım gibi doğrudan ifadeyle (değişkenin gözlenebilir olgular sayesinde doğrudan ölçülmesi anlamında değil ı ister 'yapılmaktadır' gibi dolaylı ifadeyle alınıyor olsun tüm psikolojik içerikli ölçeklerde ölçülenler tutum genel kalıbı içindedir. Tutumlarda, verilen cevapların hiçbiri için doğru-yanlış söz konusu değildir. Verilen cevap ne olursa olsun o kendine göre algılamasını ifade etmektedir, tutumunu yansıtmaktadır. Tutum ölçerlerden elde edilen ölçümler arası farklar: gerçek farkları ve sadece random hataları (burada farkı doğuran araç kökenli tüm etkenler random nitelikli olduğundan, yanlış hata olmadığından) içerdiğinden geçerlik ve güvenilirlik çalışması yapılmalıdır ve anlamlıdır. Tutum ölçeklerinde ölçek düzeyinde, faktör düzeyinde veya madde düzeyinde bir puan elde edildiği için de istatistik işlemlerin yapılması için gerekli nitelikte veri elde edilebildiğinden istatistiksel geçerlik ve güvenilirlik çalışması yapılabilir'.

Yine yargısal verilerin görüş, düşünce ve durum bildirim şeklinde alındığı anketlerde elde edilen veriler tutum ölçeklerinde veya başarı testlerinde olduğu gibi puan cinsinden değil de genellikle sınıflama türündendir. Elde edilen sayısal verilerin de bir bütün olarak anlamı yoktur. Yaşınız sorusuna verilen ' :i:'J ile aylık geliriniz sorusuna verilen 20 milyon arasında (iletilen değişken bakımdan birliktelikten bahsedilmez. Aynı değişkeni ölçmediği açıkça belli olan soruların,

örneğin birbiriyle iç tutarlığı olup olmadığını belirlemeyi düşünmek ne anlama gelir.

Güvenirliği kontrol etmek için uygulamada sıkça yapıldığı gibi aynı soruyu değişik ifadelerle anketin veya tutum ölçerin farklı verinde kontrol sorusu adı altında tekrarlama bekleme amaca hizmet eder mi? Hayır. Çünkü: eğer aynı içerik farklı ifadelerle sunulmuşsa artık farklı sorudur. Benzer' olmaları eş' olmaları anlamına gelmeyebileceği için ortaya çıkan ölçüm farklılığının ifadedeki ölçüm farklılaşmadan mı yoksa başka etkenlerden mi olduğunu ayırt edecek bir teknik bulunmamaktadır.

Bir ölçme aracı olması nedeniyle tabii ki anketlerin de geçerliği-güvenirliği söz konusudur. Ama sorun 'nasıl?' belirlenebileceğidir. Öncelikle, anketten elde edilen sınıflama türündeki verilerin bilinen güvenilirlik-geçerlik istatistik yöntemlerine uygun olmadığı dikkate alınarak 'anketlerin geçerliği ve güvenilirliği hesaplanamaz' değerlendirilmesi doğru olabilir.

Verilen cevapların doğru-yanlılığı (gerçeğe, olgulara uygunluğu veya aykırılığı) söz konusudur. Yaşınız?' sorusuna 30 yaşında olduğu halde **15** cevabını vermesi olguyu yansıtmaması bakımından yanlış bir cevaptır. Böyle yanlış cevapların alınabildiği anketten elde edilen ölçüm sonuçlarının içine yanlış hata karışmış olmaktadır. Ölçülen nesnelere arasındaki ölçüm farkları sadece ölçülenden kaynaklanan fark olmaktan çıkıp yanlış hatayı da içerecektir. Oysa istatistik ister sabit, ister değişken olsun yanlış hataları hesaplayamaz. İstatistik yanlış hatalardan arınmış, random hataları (belki) içeren verilerle çalışır. Diğer taraftan güvenilirlik, random hatalardan arınmış olmayı içeren bir kavramdır. Ölçümlerinde sürekli 1 kilo fazla tartan bir basküle çıktığımızda her seferinde sizin ağırlığımızı 60 kilo olarak belirlemesi, bu ölçümleri arasındaki tutardık onun güvenilir olduğunu göstermez. 1 kiloluk yanlış hatayı dikkate almamış olursunuz.

1.1. Geçerlik

Geçerlik, bir ölçme aracının ölçmeyi amaçladığı özelliği, başka herhangi bir özellikte karıştırmadan, doğru ve tam olarak ölçebilmesidir. Geçerlik, ölçülmek istenen değişkenin ölçülebilmiş olma derecesidir. Geçerliğin yüksek olması, ölçülmek istenen kavramın gözlenebilir nitelikteki değişkenlerle ifade edilebilmesine bağlıdır. Bu nedenle doğrudan ölçmelerde geçerlik, dolaylı ölçmelere göre daha yüksektir.

Bir ölçmenin geçerliği; objektiflik, ayırt edicilik (maddelerin ayırt etme gücü güvenilirlik değil bir geçerlik sorunudur,!, kapsamlılık, kolay uygulanabilirlik ve puanlanabilirlikten etkilenir.

Güvenirlik ve geçerlik birbirinden tamamen ayrı iki kavram olup elde edilen istatistik değerleri arasında da hiçbir ilişki yoktur. Elde edilmişleri birbirlerinden azadedir. Bir ölçme aracının mutlaka hem geçerli hem güvenilir olması anlamında birbirlerini bütünlerler'. Bunlardan biri diğerinin önkoşulu veya ikamesi değildir. Her ikisinin birlikte olmasıyla ölçme anlam bulur. Örneğin, ölçme güvenilir olması onun geçerliğini garantilemez.

Geçerlik konusu "içerik geçerliği ve yapı geçerliği" olmak üzere iki başlık altında incelenebilir.

1.1.1. İçerik Geçerliği (Content Validity)

Ölçme aracının tanımlanan davranış tepki evrenini yeterince temsil edebilmesidir. Önce ölçülmesi istenen kavram ile ilgili davranışlar evreninin çözümlenmesi yapılarak içerdiği etkenlerin açık olarak saptanması gerekir. Sonra da, her etkenin davranış evrenindeki oranına göre ölçme aracında temsil edilip edilmediğine bakılır. İçerik geçerliği ölçülen konudaki tüm boyutlardaki olası tüm maddelerden oluşan tepki (soru, madde) evrenini, ölçeğin (sınırlı sayıda madde içermesi nedeniyle bu evrenden bir örnekleme niteliğindedir! Temsil etme gücüdür. Temsil ediciliği sağlamak için her alt boyutu içeren maddeler ölçekte yer almalıdır ve tepki evrenindeki oranı veya önem ağırlığı ölçeğe de yansmalıdır.

Örneğin, matematik bilgisi ölçülmek istendiğinde, öğrenciye sadece toplama işlemini sormak onun dört işlemin diğer öğelerini bilip bilmediğini anlamaya yeterli olmaz. Bu nedenle ölçülmek istenen amacın tümünü kapsayan sorular sorulmalıdır.

İçerik geçerliğini, aynı değişkenle ilgili bir başka ölçeğin sonuçları ile korelasyon arayarak belirlemek, sonuçta istatistik bir değere dayanarak içerik geçerliğinden bahsetmek yaygın bir kanı olmakla beraber yanıltıcıdır. Her şeyden önce 'ölçüt olarak alınan ölçeğin içerik geçerliği nasıl hesaplanmıştır?' sorusuna alınacak 'bir başka ölçme aracını ölçüt olarak' cevabı kısır döngüyü gözler önüne serip içerik geçerliğini hesaplama hevesini boşa çıkaracaktır.

1.1.1.1. Yüzeysel Geçerlik (Face Validity)

Ölçme aracının hangi değişkeni ölçtüğü hakkındaki uzman görüşüdür. Geçerlik seviyesini sayısal değerle belirtme olanağı yoktur. Sadece kanaatlere göre bir kabul söz konusudur. Geliştiren ölçme araçlarında ilk başvuru olan geçerlik türüdür. Bir veya birkaç uzmanın görüşüne başvurulurken, ölçme aracının kullanılacağı amaç için gerekli veriyi toplayacak durumda olup olmadığı yönünde alınan bilgi yüzeysel geçerlik için yeterli görülmektedir. Yüzeysel geçerlik ölçme aracının hangi değişkeni ölçtüğünü değil, ölçer gibi göründüğünü belirler.

1.1.1.2. Uygulama Geçerliği (Empirical Validity)

Ölçülmeye çalışılan değişkenin gerçek hayattaki gözlemlenebilir belirtileri ile ölçme sonuçları arasındaki uyum. uygulama geçerliğidir. Uygulama geçerliğini belirlerken dış ölçüt ile ilişkisine bakılır. Ölçüt, ölçmek istenen davranışı ne dereceye kadar yansıtırsa, bulunan ilişki ölçme aracının geçerliği hakkında o derecede sağlam bilgi verir. Ölçüt saptanması ölçme aracının niteliğine ve ölçmenin amacına göre farklılık gösterir. Örneğin okulda alınan notlar ile iş hayatındaki başarı arasında yüksek bir ilişki bulunmuşsa okuldaki notların uygulama geçerliği olduğu söylenebilir. Burada ölçüt, iş hayatındaki başarıdır. Atılganlık ölçeğinde yüksek puan alan bireyler, günlük yaşamlarında yakın

çevrelerince pısrık olarak tanınıyorlarsa atılğanlığı ölçmek için kullanılan ölçğin pratik geçerliğinin olduđu düşünülemez.

Ölçeklerde uygulama geçerliği hakkında bazı bilgiler elde etmek için su yaklaşımlar kullanılabilir.

- Ölçülen değışken bakımından bireyleri bu açıdan tanıdıklarına inanılan başka kişiler, onları sıralar veya sınıflar (atılğan-çekingen vb). Ölçülen bireyler hakkında gözlem yoluyla da aynı veriler toplanabilir. Elde edilen bu veriler ile ölçüm sonuçları arasındaki yüksek uyum kullanılan ölçğin uygulama geçerliğinin olduğunu ifade edebilir,

- Ölçülen değışken bakımından uç değerlere sahip oldukları hakkında yaygın kanaatin bulunduğu birey ya da gruplara geliştirilen ölçek uygulanır. Hakim olan kanaat ile ölçek sonuçları arasındaki uyum yüksekse uygulama geçerliğinin olduğu, ölçğin uygun içerikte olduğu yönünde ipucu elde edilmiş olur. Örneğin akademik yetenekler ölçümünde kullanılacak bir test hazırlandığında araç, zihinsel gücü bakımından yukarı (üstün okulları) ve aşağı (alt özel sınıflar, vb) seviyeleri temsil eden bireylere uygulanır. Yukarı seviye diye bilinenler ölçekte de üst puanları almış, aşağı seviye diye bilinenler de ölçekte alt puanları almış ise, ayrıca bu iki karşıt grubun ölçek puanları karşılaştırıldığında üstünler-yönünde manidar bir farklılık belirmişse ölçğin uygulama geçerliği olduğu yönünde veri elde edilmiş olur.

Uygulama geçerliği ilgilendiği zaman dilimine göre "halihazır geçerlik" veya "kestirirsel geçerlik" adını alır. Şimdiki zamandaki uygulama geçerliği araştırılıyorsa hâlihazır geçerlik, gelecek zamandaki uygulama geçerliği araştırılıyorsa kestirirsel geçerlik söz konusudur.

Bu uygulama geçerliği tür ve tekniklerinden hangisiyle olursa olsun elde edilen geçerlik, içerik geçerliği hakkında fikir verecektir. Uygulama geçerliği içerik geçerliğinin dışında bir anlam taşımaz Çünkü uygulama geçerliğinde kullanılan ölçütler ölçülen değışkenle yani içerikle ilgili olmak zorundadır. Uygulama geçerliğini belirlemede biricik sorun objektif veriler- elde ettiğimiz (ya da etmeye çalıştığımız) ölçek sonuçlarının uygulamayla karşılaştırmasının yaygın kanı. Ön kabul gibi sübjektif verilere dayandırılmasıdır.

1.1.1.2.1. Kestirirsel Geçerlik (Predictive Validity)

Değişkenler arası ilişkileri anlayarak olayları önceden tahmin edebilip kontrol altına almak bilimin nihai amacıdır. Bilimin ortak amaçlarından olan geleceğe yönelik tahmin, ailem araçları için de söz konusudur-. Sosyal bilimlerde ölçme araçları, genellikle bireylerin gelecekte gösterecekleri davranışları tahmin etmek amacıyla geliştirilirler. Bir zeka testi ilerideki akademik başarıyı, bir yetenek testi meslekteki başarıyı, bir politik tutum ölçüğü belli niteliklerdeki kişilerin hangi partiye oy vereceklerini önceden tahmin etmeye olanak sağlayabilmelidir.

Kestirirsel geçerliğin saptanmasında izlenen yol şöyledir. Ölçme aracı uygulanarak sonuçlar alınır. Ölçülen niteliğin belirgin olarak görülebileceği yeterli bir süre beklenir. Belirlenen ölçüt açısından da uygulamadaki durumun

değerlendirmeleri yapılır. Geliştirilen testin sonuçları ile ölçüt değişken sonuçları arasında manidar bir ilişki bulunmuşsa testin tahmin geçerliği olduğu kabul edilir.

1.1.1.2.2. Halihazır Geçerlik (Concurrent Validity)

Her şeyden önce belli bir zaman aralığı gerektirdiğinden çoğu kere önce ölçme aracının uygulanması, bir süre sonra ölçüt değişken değerlendirmesinin yapılması mümkün veya uygun olmayabilir. Bu durumlarda halihazır geçerlik ile yetinilir. Halihazır geçerlikte, ölçme aracı kullanıldığında, ölçüt değişken değerlendirmesi de aynı zamanda yapılır.

1.1.1.3 Ölçeğin İctutarlığı

Ölçeğin içerik geçerliği hakkındaki değerlendirmelerde enendi ölçütlerden biri de aracın içtutarlığıdır. İctutarlık temelde bir güvenilirlik sorunudur. Ölçkleme varsanımlarından tok boy ut tuluğu il karşılama durumunu belirlemek için hesaplanır. Güvenirlik başlığı altında ayrıntıları verildiği gibi bir faktörü oluşturan maddeler birbiriyle ne kadar yüksek ilişki gösterirse tek boyutlu oldukları, yani ölçülen değişken her ne ise de aynı değişkeni ölçtükleri kararına varılır.

İlişki ne kadar yüksekse o ölçüde tek boyutluluk hakimdir. Güvenirlik açısından durum böyle olmakla beraber geçerlik açınsındansa tepki evreni temsil etmek bakımından maddeler birbirinden ne kadar farklı ise aynı değişkenle ilgili o kadar çok örnekleme yapılmıştır. Ölçeği oluşturan faktörler veya faktörü oluşturan maddeler hem birbiriyle ilişki olacak (ıctutarlık için) hem de birbirinden farklı olacak (içerik geçerliği için). Çünkü birbirinden farklı olmayan maddelerin ölçekte bulunmasının cevaplayanı bunaltmaktan başka bir etkisi olmayacaktır, faydası ise hiç. İlk bakışta çelişki gibi gözükten bu durum istatistik teknikler açısından çok normaldir. İctutarlık için ilişki belirleyen korelasyon tekniklerinin, içerik geçerliği için ise fark belirleyen (t-testi, varyans analizi) tekniklerinin her ikisinin de kullanılması gereğini hatırlatmaktadır. Çünkü istatistiksel olarak iki veri dizisi birbiriyle hem ilişkili olabilir ama aynı zamanda da farklı olabilir

İctutarlık için beliren bu değerlendirmeye 'yapı geçerliği" başlığında açıklanan faktör analizi tekniği kısmen (çünkü tüm korelasyon-kovaryanslara göredir) kolaylık getirmektedir. Faktör analizinin hesaplarına tarzı olarak birbiriyle en fazla ilgili maddeler aynı faktör etrafında toplanarak ilişkiyi, birden fazla faktör olması ise farkı yansıtmaktadır.

1.1.2. Yapı Geçerliği (Construct Validity)

Soyut kavramlara yönelik ölçmelerde önce ölçülen kavramı tanımlayan kuramlardan biri tercih edilir. Böylece ölçülmek istenen kavramın yapısı belirlenir. Bu kuramsal yapıya göre gözlenebilir değişkenler ortaya konur. Son olarak bu gözlenebilir değişkenlerini ifadelendiren maddeler yazılarak ölçek hazırlanır. Ölçek geliştirildiğinde maddelerin hangi faktörleri temsilen

yazıldığını arařtırmacı bilmektedir. Yani teorik yapıya baęlı ölçek yapısı belirlidir.

Kavram ile deęişkenin özdeş olmadığı (ki sosyal bilimlerdeki ölçmelerin tümünde kavram ve deęişken özdeş deęildir), bunun yanında ölçmenin gözlenebilir deęişken deęil de onun belirtileri üzerinden yapıldığı durumlarda, ölçülenin ölçülmek istenen kavramın belirtileri olup olmadığının belirlenmesi gerekir. Yapı geçerlięi, bu belirleme çalışmalarının sonucudur Dięer taraftan yapı geçerlięi, ölçeęin temelindeki kuramların geçerlięi ile de ilgilidir.

Ölçeęin uygulanmasından elde edilen veriler, "faktör analizi" istatistik teknięi ile işlenir. Çeşitli faktör analizi teknikleri olmakla beraber (varimax rotated faktör analizi en yaygın kullanım bulanıdır) temel mantık aynıdır. Maddeler arasındaki olası tüm ilişkiler hesaplanıp bir korelasyon matrisi oluşturulduktan sonra aralarında beliren ilişki yapılarına göre farklılık gösterdikleri oranda faktör adedi belirir, her maddenin her faktör içindeki aęırlığını ifade eden katsayılarından (maddenin faktörle ilişkisi olarak yorumlanabilir) oluşan bir tablo elde edilir. Bu tablo üzerinde yapılan deęerlendirmede her madde en yüksek aęırlığı hangi faktörde bulmuşsa o faktörün kapsamında olmasına karar verilir. Bu maddenin her-hangi bir faktöre girebilmesi için ulaşması gereken asgari bir deęer istatistiksel olarak söz-konusu deęildir, konulan barajlar keyfidir. Sanıldığı gibi aksine faktör analizi sonunda řu maddeler řu faktörü oluşturur gibi kesin bir yargıyı gerektirecek istatistik karar çıkmaz. Faktör-madde ilişkisi pratikte, ideal deęerini bulmadığından bir madde birden fazla faktörle yüksek ilişki gösterebilir. Bazı yabancı test uygulamalarında görüldüğü gibi bir madde birden fazla faktörde yer alabilir. Bu durumun uygulamada bazı deęerlendirme ve yorumlama güçlüklerine yol açabileceğini de göz önünde tutmak gerekir.

Ölçeęin başlangıçta belirlenen teorik yapısı ile uygulama sonuçlarına dayanan faktör analizi işleminin sonrasında beliren pratik yapı birbirine uyumlu ise ölçeęin yapı geçerlięi vardır. Hem ölçeęin dayandığı kuramsal yapı doğrudur hem de ölçek bu kuramsal yapıyı ölçebilecek niteliktedir.. Teorik yapı ile pratik yapıyı karşılařtırmaya, ikisi arasındaki uyumu belirlemeye hizmet eden bir istatistik teknik henüz maalesef yok.

Öğrencilerin ders çalışma alışkanlıklarını belirlemeye çalışan arařtırmacının "zamanı kullanma" boyutu ile ilgili olduğunu düşünüp yazdığı maddeler faktör analizi sonuçlarında da bir faktör olarak belirmişse yapı geçerlięinin olduğu ve boyutun "zamanı kullanma" kavramını ölçtüğü düşünülebilir. Ama bazı maddeler başlangıçta düşünülenin dışında maddelerle bira raya gelmişse yapı geçerlięi yoktur. Böyle bir durumda ya ölçek, önerilen kuramsal yapıyı ölçmekten uzaktır ve/veya ölçeęin dayandırıldığı kuramsal yapı yanlıştır.

Ölçeęin yapı geçerlięi için faktör analizinin yapılması şarttır. Bunun yanında sonuçları desteklemek amacıyla, daha önceden yapı geçerlięi faktör analiziyle belirlenmiş ve aynı kuramsal yapıya dayanan bir başka ölçme aracı ile birlikte uygulanarak sonuçlar arasındaki uyum deęerlendirilebilir

Faktör analizi ile ilgili iki yanlış kanı vardır. Birincisi; herhangi bir kuramsal yapı belirlemeden faktör analizi sonuçlarını olduğu gibi kabul edip ölçeği yapılandırmak. Bu elbette kolaylık sağlamaktadır, ama başlangıçta belirlenen bir kuramsal yapı olmadığından bu tarz hareketin yapı geçerliğini sağladığını düşünmek yanlış olur. herhangi bir yapı belirlenmemiştir ki geçerli olup olmadığına karar verilmiş olsun. İkincisi; yabancı dilde ve ülkede geliştirilen bir ölçek Türkçeye çevrilirken (dilsel adaptasyon,) bu ölçeği oluşturan maddelerin içeriği aynı kalırken faktör analizi yapmakla yeniden yapılandırarak, yapısal adaptasyon hakkını da kendinde görmektir. Öncelikle, ölçeği geliştiren biz olduğumuz halde geliştirenin belirlediği yapıyı değiştirmek, bozmak bilimsel etik bakımından ne ölçüde uygundur? Çünkü ölçeğin yapılandırılması belli bir kuramsal yapıyı yansıtır, başkası adına bu kuramsal yapıyı biz nasıl belirleyebiliriz. Adapte eden biz olsak bile ölçeğin yasal ve bilimsel sahibi orijinalini geliştirendir. Var olan orijinal kuramsal yapı eğer beğenilmiyorsa dilsel uygulama çalışmasına girmek gereksiz, yok eğer beğeniliyorsa niye yapı değiştirilmek istensin. Faktör analizi ölçek geliştirilmesinde kullanılabilecek bir tekniktir, adaptasyonunda değil. Uyarlaması yapılan ölçeğin orijinal ile eşdeğer olduğunu iddia edebilmek için gerekli koşullar eşdeğerlik başlığı altında incelenmiştir Ölçeğin yapısının değiştirilmesi eşdeğerliği ortadan kaldırır, genel hatlarıyla aynı değişkeni ölçen farklı bir ölçek ortaya çıkmış olur. Dolayısıyla bu iki ölçekten elde edilen sonuçlar hiçbir şekilde karşılaştırılmaz. Bu ise bilimin evrenselleşmesi doğrultusundaki kapsayıcı genellemelere ulaşılmasına bir engeldir.

1.1.3. Ayırt etme Gücü

Bir maddenin ayırt etme gücü, ölçülen değişken bakımından birimler arası farklılığı ne ölçüde ortaya çıkarabildiği ile ilgilidir. Ölçmenin temel amacı ölçülen nesnelardaki farkı yakalayabilmek olduğuna göre ayırt etme gücü ayrı bir önem kazanmaktadır. Bir maddeye, herkes aynı cevabı vermiş ise diğer özellikleri ve önemi ne olursa olsun kimseyi diğerinden ayırt etmediği için maddeyi ölçekte tutmanın bir yararı yoktur. Ayırt etme gücü zayıf maddelerin ayıklanması ile ölçek daha kısa ama etkili bir hale getirilmiş olur.

Maddelerin ayırt etme gücünün analizi için bireylerin, ölçekten aldıkları toplam puanı belirlenir ve bu toplam puana göre en büyükten en küçüğe doğru sıralanırlar, O maddeyi cevaplayan bireylerin %27 sinin kaç kişi olduğu belirlenir. Sıralamanın en üstündeki %27'lik grup (n_u) ile en altındaki % 27'lik grup (n_a) belirlenir. Bu işlem öncesi hazırlıktan sonra ölçek değer ayrımı veya t-testi ile gereken sınama yapılır. İstatistik yetkinliği bakımından elbette t-testi tercih edilir. Ayırt etme gücünü belirleme tekniklerindeki temel yaklaşım, testin toplamında yüksek puan alanların incelenen madde de yüksek puan almaları gereğinin karşılanıp karşılanmadığıdır. Aynı şekilde ölçeğin toplamında düşük puan alanlar grubunda yer alan bireylerin madde puanlarının da düşük olması gerekir. Üst ve alt çeyrekte yer alan bireylerin her madde için bu konumlarını koruyup korumadıkları test edilir. Başarı testlerinden örneklemek gerekirse çok kolay olduğu için hemen herkes tarafından doğru cevaplandırılan veya çok zor

olduğu için hemen hiç kimse tarafından doğru cevaplandırılmayan sorular ayırt etme gücü bakımından zayıftır.

Tıpkı ölçeğin içtutarlılığındaki gibi ayırt etme gücü için de geçerlik-güvenirlik ikilemi vardır İçtutarlık maddeler arası, ayırt etme gücü ise madde içinde bireyler arası değişkenlik ile ilgilidir. Geçerlik açısından ayırt etme gücünün yüksek olması gerekir. Yani bireyleri puanlama bakımından birbirinden ne ölçüde farklılaştırabiliyorsa o ölçüde ayırt edicidir. Bu ise büyük ölçüde istatistiksel değer olarak varyansın büyümesine yol açar. Diğer taraftan ölçmenin standart hatası konusunda görüldü ki, varyans büyüdükçe standart hata da büyümektedir. Yani güvenilirlik ölçütlerinden birine göre güvenilirlik düşmektedir. Hatanın, güvenilirliği yükselttiğini düşünmek mümkün mü? Bazı kaynaklarda rastlandığı gibi ayırt etme gücü artıkça güvenilirlik artmaz, zincirleme olarak standart sapma büyür, standart bata büyür ve güven aralığının sınırları genişler. Güven aralığının sınırlarının genişlemesi güvenirlığın arttığına delalet etmez, sadece gerçek ölçüm sonucunu tahmin gücünün zayıfladığının ifadesidir.

TABLO 1. AYIRT ETME GÜCÜ HESAPLANIŞI (1. ÖRNEK)

Öğrenci	2,soru	Ölçek Toplamı	2. sorunun istatistik değerleri				α	α^2
			n	$\sum x$	$\sum x^2$	μ		
A	3	40 ü	3	10	34	3.33	0.577	0.333
B	4	40 ü						
C	3	35 ü						
D	2	33						
E	4	30						
F	3	29						
G	4	28						
H	2	26 a	3	4	6	1.33	0.577	0.333
I	1	16 a	alt çeyreklik					
İ	1	16 a						
Toplam			10	27	85	2,70	1.160	1.344

1.1.3.1. Ölçek Değer Ayrımı (D_a)

Ölçek değer ayrımı ölçek toplam puanına göre üst çeyrekte yer alanların aritmetik ortalamasından alt çeyrekte yer alanların aritmetik ortalamasının çıkarılmasıyla elde edilen katsayıdır.

$\mu_{\bar{u}}$: üst çeyreklikte yer alanların aritmetik ortalaması

μ_a : alt çeyreklikte yer alanların aritmetik ortalaması

D_a : maddenin ölçek değer ayrım katsayısı

$$D_a = \mu_{\bar{u}} - \mu_a \quad [11]$$

Tablo 1'deki örneğe göre 10 kişilik bir grubun üst çeyreğinde yer alan 3 kişinin incelenen 2. maddeden aldıkları puanların aritmetik ortalaması 3.3 (10/3), alt grubun 1.33 (4/3)'dür. Bu maddenin ölçek değer ayrımı $D_g = 3.33 - 1.33 = 2$ 'dir.

Aynı çözümlene tekniğini bir başka form içinde örnekleyen Tablo 2'deki verilere göre 152 kişilik bir gruba uygulanan bir ölçekteki bir maddenin ölçek değer ayrımı 0.64 bulunmuştur.

Ölçek değer ayırım katsayısı ne kadar büyükse ilgili madde o ölçüde ayırıcı değere sahiptir. Ancak bununla ilgili istatistiksel bir manidarlık sınaması söz konusu değildir. Ne kadarlık miktardaki ölçek değer ayırımı, ayırt ediciliğin var olduğunu ifade ederin kararı objektif ölçütlere bağlanmamıştır. Ayrıca dizilerde sadece merkezi eğilim değerini dikkate alıp dağılım değerleriyle ilgilenmemesi nedeniyle grupları karşılaştırmada istatistiksel yetkinliği yüksek değildir. Bu nedenle aynı amaca hizmet eden ve istatistik açıdan manidarlık sınamasının yapıldığı t-testini, ayırt ediciliği belirlemede kullanmak çok daha yararlı olacaktır.

TABLO 2. AYIRT ETME GÜCÜ HESAPLANIŞI (2.ÖRNEK)

	Ağırlık	Üst Çeyreklik			Alt Çeyreklik		
	x	f	fx	fx ²	f	fx	fx ²
Kesinlikle katılıyorum	5	5	25	125	1	5	25
Katılıyorum	4	10	40	160	5	20	80
Kararsızım	3	20	60	180	20	60	180
Karşıyım	2	5	10	20	10	20	40
Kesinlikle Karşıyım	1	1	1	1	5	5	5
Toplam		n	Σfx	Σfx²	n	Σfx	Σfx²
		41	136	486	41	110	330
$\mu = \bar{x} = \frac{\Sigma X}{n}$		$\mu_{\bar{u}} = \frac{136}{41} = 3.32$			$\mu_{\bar{a}} = \frac{110}{41} = 2.68$		
$\sigma^2 = \frac{n \Sigma X^2 - (\Sigma X)^2}{n(n-1)}$		$\sigma_{\bar{u}}^2 = \frac{41 \times 486 - 136^2}{41 \times 41} = 0.851$			$\sigma_{\bar{a}}^2 = \frac{41 \times 330 - 110^2}{41 \times 41} = 0.851$		
n ≥ 30 ise (n-1) yerine (n)							
$\Sigma fx = \Sigma x$ $\Sigma fx^2 = \Sigma x^2$	Madde Ölçek Değer Ayırımı: 3.32-2.68=0.64						

1.1.3.2. t-testi (Kritik Oran)

İstatistikteki en yaygın tekniklerden olan t-testinin, karşılaştırılan grupların bazı özelliklerine göre farklı modelleri vardır. Maddelerin ayırt ediciliğini belirlemede kullanılacak olan t-testleri ilişkisiz gruplar için kullanılmaktadır. T-testi genel form u şöyledir.

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sigma_F} \quad [1.2]$$

Farkın standart hatası (σ_f) örneklem büyüklüğüne bağlı olarak farklı hesaplanmaktadır. Eğer $n \geq 30$ ise

$$\sigma_f = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \quad [1.3]$$

Eğer $n < 30$ ise

$$\sigma_f = \sqrt{\frac{\sigma_1^2}{n_1 - 1} + \frac{\sigma_2^2}{n_2 - 1}} \quad [1.4]$$

t-testinin bu şekilde hesaplanmasında aritmetik ortalama ve standart sapma hesaplamalarında kesirli kesimdeki yuvarlamalardan ötürü sonuçta küçük farklılıklar görülebilir. Daha duyarlı bir (t) değerine ulaşılmak isteniyorsa formül [1.5] kullanılabilir. Eğer büyük örneklere göre işlem yapılacaksa formüldeki (n_1-1) ve (n_2-1) yerine (n_1) ve (n_2) kullanılacaktır.

$$t = \frac{\sum X_1 - \sum X_2}{\sqrt{\frac{n_1 \sum X_1^2 - (\sum X_1)^2}{n_1 - 1} + \frac{n_2 \sum X_2^2 - (\sum X_2)^2}{n_2 - 1}}} \quad [1.5]$$

Tablo 1'deki veriler formül [1.4]'e uyarlandığında $\sigma_f = \sqrt{\frac{0.333}{3-1} + \frac{0.333}{3-1}} = 0.577$ bulur.

nur. İşleme formül [1.2] uygulanarak devam edilince $t = \frac{3.33 - 1.33}{0.577} = 3.46$ bulunur. $Sd = (n_1 - 1) + (n_2 - 1) = 4$ için bir yönlü olarak manidarlık sınaması yapıldığında $3.747 > 3.46 > 2.132$ olduğundan sözkonusu maddenin 0.95 güven düzeyiyle üst ve alt çeyreklikleri ayırt ettiği, t katsayısının pozitif işaretli olması nedeniyle bu farklılığın beklenen yönde olduğu, bütün bu istatistik analizler sonucu maddenin ayırt edici özelliği olduğu ve bu ölçüte göre ölçekte yer alabileceği kararına varılır. Manidar çıkmayan maddeler ölçekten kesinlikle çıkarılır. Eğer ölçmeye çalıştığı ayrıntıyı ölçen başka maddeler ölçekte bulunmuyorsa aynı amaca yönelik başka ifadeler, başka maddeler hazırlanıp, ölçek yeniden yapılandırılarak geçerlik-güvenirlik işlemleri yenilenir. Geçerlik ve güvenirlik işlemleri, bu açılardan aranan ölçütler karşılancaya kadar devam ettirilir. İşlemler sonucu ölçekten bazı maddelerin atılması ile ölçek yeni bir yapıya kavuşmuştur. Bizim yaptığımız geçerlik-güvenirlik işlemleri eski ölçek için geçerliydi, ölçeğin işlemler sonucu aldığı son yapı içindeki geçerlik-güvenirliği nedir? sorusunu araştırmacı olarak merak etmez misiniz? t testine göre manidar çıkıp ayırt edici olduğuna karar verilen bir madde ölçekte mutlaka yer alacak demek değildir. Bu sadece ayırt edicilik bakımından uygun olduğunun göstergesidir, diğer ölçütler örneğin içtutarlık sonuçları olumlu çıkmamışsa ayırt edici olmasına rağmen bu madde ölçekten ayıklanacaktır. Bazen t değeri istatistiksel bakımdan manidar ve fakat negatif çıkabilir, bu alt çeyreklikte yer alanların üst çeyreklikte yer alanlardan daha yüksek puan aldıklarını gösterir, yani beklenen durumun tam tersine bir sonuçla karşılaşmıştır. Aslında sözkonusu madde alt ve üst çeyreklikleri ayırt edicidir, ama ters yönde. Bu, böyle mad-

delerin puanlandırılmasında hata olduğu şeklinde yorumlanmalıdır. Bu maddenin puanlamasının ters çevrilebilme imkanı varsa, örneğin Likert türü bir ölçek ise ve "her zaman" seçeneği 1 puan "hiçbir zaman" seçeneği 5 puan şeklinde değerlendirilmiş iken "her zaman seçeneği 5 puan "hiçbir zaman seçeneği 1 puan alacak şekilde değerlendirmenin yönü deriştirilir. Bu işlem ayırt ediciliği beklenen yöne çevirecek ama toplam puan üzerinde yapacağı etki nedeniyle diğer tüm geçerlik-güvenirlik hesaplarında da büyük bir ihtimalle değişikliğe yol açacaktır. Analizlerin bitiminden sonra geçerlik-güvenirlik çalışmalarının yinelenmesi gereği belirmektedir.

Tablo 1'deki veriler formül [1.5]'e göre çözümlendiğinde

$$t = \frac{10 - 4}{\sqrt{\frac{3 \times 34 - 10^2}{3-1} + \frac{3 \times 6 - 4^2}{3-1}}} = 4.24$$

olarak bulunur ve manidarlık tablosu değeri olan 3.747 den büyük olduğu için .99 güven düzeyiyle ayırt edicidir.

Tablo 2'deki veriler formül [1.3]'e uyarlandığında $\sigma_t = \sqrt{\frac{0.851}{41} + \frac{0.851}{41}} = 0.204$ bulunur.

Tablo 2'deki veriler formül [1.3]'e uyarlandığında $\sigma_t = \sqrt{\frac{0.851}{41} + \frac{0.851}{41}} = 0.204$ bulunur.

İşleme, formül [1.2] uygulanarak devam edilince $t = \frac{3.32 - 2.68}{0.204} = 3.137$ bulunur. $Sd = (n_1 - 1) + (n_2 - 1) = 80$ için bir yönlü olarak manidarlık sınaması yapıldığında .01 düzeyinde manidar olduğu bulunur. Bu maddenin ayırt etme gücü yüksektir kararına varılır.

Tablo 2'deki veriler formül [1.5]'e göre çözümlendiğinde $t = \frac{136 - 110}{\sqrt{\frac{41 \times 486 - 136^2}{41} + \frac{41 \times 330 - 110^2}{41}}} = 3.113$ olarak bulunur ve .99 güven düzeyiyle ayırt edicidir.

Bazı kaynaklarda kritik oran adıyla verildiği şekil Tablo 2'deki verilere uygulanırsa

$$\text{Kritik Oran} = \frac{\sum fx_1 - \sum fx_2}{\sqrt{\sum fd_1^2 + \sum fd_2^2}}$$

$$\sum fd_1^2 = \sum fx_1^2 - n_1 (\bar{X}_1)^2$$

$$\sum fd_2^2 = \sum fx_2^2 - n_2 (\bar{X}_2)^2$$

$$\text{Kritik Oran} = \frac{136 - 110}{\sqrt{(486 - (41(3.32^2))) + (330 - (41(2.68^2)))}} = 3.12$$

sonucun yine aynı olduğu kolaylıkla görülür.

Bir maddenin ayırt ediciliğini belirlemek için elde edilen t katsayısının (kritik oranın) ne kadar büyük olduğuna bakarak değil, serbestlik derecesine göre manidar olup olmamasına göre karar verilir, Manidarlık düzeyi ne kadar yüksekse ayırtıcı özelliği de o ölçüde yüksektir. Hatırlanması gereken nokta bir maddenin ayırt edici olması onun güvenilir olmasını (örneğin içtutarlığı) garantilemez.

1.2. Güvenirlilik

Bir ölçme aracıyla farklı zamanlarda elde edilen ve aynı nesnelere ilgili olan bir grup ölçümle ikinci grup ölçüm arasındaki tutarlık eğilimine o aracın güvenirliliği denir. Güvenirlilik; aynı değişkenin bağımsız ölçümleri arasındaki kararlılıktır, aynı süreçlerin izlenmesi ve aynı ölçütlerin kullanılması ile aynı sonuçların alınmasıdır, ölçmenin random hatalardan arınmış olmasıdır.

Güvenirliliğin yüksek olabilmesi, ölçmede izlenen süreçler ile kullanılan ölçütlerin ayrıntılı olarak belirlenebilmesine bağlıdır. Dolaylı ölçmelerin yaygın olarak kullanıldığı sosyal bilimlerde güvenirliliği yükseltmek için çok sayıda ölçüt

kullanılmaya çalışılır, madde veya soru sayısı arttırılır. Böylece random hataların birbirini dengelemesiyle güvenilirliği yüksek sonuçlar alınabilir.

Birbirini izleyen ölçmelerde, bireyin grup içindeki durumundaki tutarlık aranır. Bu yaklaşımda aynı nesnelere ilgili iki ölçüm arasında korelasyon hesaplanır ve bulunan korelasyon katsayısı, güvenilirlik katsayısı olarak adlandırılır.

Güvenirlik katsayısı, gerçek ölçümlerin varyansının gözlenen puanların (gerçek ölçüm ve hata) varyansına oranıdır.

σ^2_g : gerçek puanlar dağılımı varyansı

σ^2_t : gözlenen toplam puan dağılımının varyansı

σ^2_h : hata varyansı (standart hata karesi)

r_r : elde edilen puanların güvenilirlik katsayısı

$$r_r = \frac{\sigma^2_g}{\sigma^2_t} \quad [2.1]$$

$$\sigma^2_t = \sigma^2_g + \sigma^2_h$$

$\sigma^2_g = \sigma^2_t - \sigma^2_h$ olacağından Formül 1,

$$r_r = \frac{\sigma^2_t - \sigma^2_h}{\sigma^2_t} \quad [2.2]$$

halini alır.

Hatasız gerçek ölçümlerin bilinmesi mümkün olmadığından güvenilirlik dolaylı yoldan belirlenmeye, tahmin edilmeye çalışır. Güvenirlik tahmininde kullanılan iki yaklaşım vardır: Ölçümler arasındaki ilişkiyle ilgilenen **güvenirlik katsayıları** ve ölçme hatalarının büyüklüğüyle ilgilenen **ölçmenin standart hatası** güvenilirlik konusundaki iki temel yaklaşımı oluşturur. Güvenirlik katsayıları ölçülen bilimin grup içindeki konumundaki kararlılığı-136 nı dikkate alır. Ölçmenin standart hatası ise ayrıca; ya farklı birimlerin ölçümleri yada aynı birimin farklı ölçümleri arasındaki değişkenliği (kararlılığı) gösterir.

1.2.1. Ölçmenin Standart Hatası

Güvenirlik katsayısının hesaplanabilmesi için; ölçüm sonuçlarının ne kadarının gerçek varyansı ne kadarının hata varyansı olduğunu bulabilmek, bunun için ise gözlenen değerle birlikte gerçek değeri de bilmek gerekir. Bu söz konusu olamayacağına göre hatanın büyüklüğü ve güvenilirlik katsayısı bazı tekniklerle

tahmin edilmeye çalışır. Random hata miktarının tahmini için aynı birey hakkında aynı ölçme aracından elde edilmiş birden fazla veriye veya aynı ölçme aracından birden fazla birey için elde edilmiş veriye gerek vardır. Bu anlamda, elde edilen ölçümler birbirinden farklı olabilir ve bunlar belli bir dağılım gösterir. Standart sapma olarak ifade edilen, ölçümlerdeki bu değişme, gerçek hata payının standart hata olarak hesaplanmasında kullanılır. Puanların gösterdiği dağılımın standart sapmasından, puanların ortalamasının standart hatası hesaplanabilir. Standart hata küçüldükçe ölçmenin güvenilirliği artmış olacaktır. Ölçmenin standart hatası, özellikle, bir testteki çeşitli puanların ve puanlar arasındaki farkların güvenilirliği konusunda verilecek kararlarda kullanılır.

Ölçmenin standart hatası formülü şöyledir.

$$\sigma_m = [8\sqrt{(1-0.92)}] = 2.26 \text{ olacaktır.}$$

$$z = 2.58 \text{ (0.99 güven düzeyi için normal dağılım puanı)}$$

$$x = 60 \text{ (bireyin aldığı puan)}$$

değerleri Formül [2.4]'e yerleştirildiğinde

$$R = 60 \pm (2.58 * 2.26) = 50 \pm 5.83$$

σ_m : ölçmenin standart hatası

σ : elde edilen puanların standart sapması

r_r : elde edilen puanların güvenilirlik katsayısı

$$\sigma_m = \sigma\sqrt{(1-r)} \quad [2.3]$$

Standart hata gerçek ölçümlerin belli olasılıklarla içinde olabileceği puan aralığının belirlenmesinde kullanılır.

σ_m = : ölçmenin standart hatası

z : normal dağılım değeri

X : elde edilen puan

R : elde edilen puanın aralığı

$$R = X \pm (z * \sigma_m) \quad [2.4]$$

Test puanının güvenilirliğini belirlemeyi bir örnekle açıklama için;

$\sigma = 8$ (elde edilen puanların standart sapması)

$r_r = 0.92$ (elde edilen puanların güvenilirlik katsayısı)

değerleri Formül [2.3]'e yerleştirildiğinde

$$\sigma_m = [8\sqrt{(1-0.92)}] = 2.26 \text{ olacaktır.}$$

$$z = 2.58 \text{ (0.99 güven düzeyi için normal dağılım puanı)}$$

$$x = 60 \text{ (bireyin aldığı puan)}$$

değerleri Formül [2.4]'e yerleştirildiğinde

$$R = 60 \pm (2.58 * 2.26) = 50 \pm 5.83$$

olacaktır. Diğer bir ifadeyle bu ölçüm sonuçlarına göre, bireyin gerçek puanı, 9c99 olasılıkla 54,17 ve 65.83 arasındadır. Güven aralığının sınırları birbirine ne kadar yaklaşırsa pratik değeri o ölçüde artar. Bir ölçüm dizisinin standart hatası, elde edilen varyansın hangi oranda bireyler arası gerçek farklılığı yansıttığını gösterir.

Güvenirlilik katsayısı ve ölçmenin standart hatası ilişkisini görmek için yukarıdaki veriler Formül [2,2]'ye yerleştirilince

$$r_r = \frac{8^2 - 2.26^2}{8^2} = 0.92$$

çıkacaktır. Yani ölçmenin standart hatası ölçümdeki random hata miktarını belirlemektedir ve güvenirlilik katsayısının yükselmesi standart hatanın, ölçümlerin varyansı içindeki payının küçülmesine bağlıdır. Dolayısıyla standart hata ne kadar küçük ise ölçmenin güvenirliliği de o ölçüde artıyor demektir. Örnekte de görüldüğü gibi güvenirlilik katsayısı ve standart hata, güvenirlilik konusunda birbirini destekleyen iki ayrı sonuçtur.

1.2.2. Güvenirlilik Katsayıları

Güvenirlilik katsayıları genellikle korelasyon teknikleri ile hesaplanır. Korelasyon katsayısı -1.00 ile + 1.00 arasında değişmekle birlikte, güvenirlilik katsayılarının pozitif değerli olması beklenir, negatif değerli katsayı güvenirsizliğin bir göstergesidir. Değer +1'e yaklaştıkça güvenirliliğin yüksek olduğu kabul edilir. Güvenirlilik katsayıları tekniklerinin mani-darlıkları ile ilgili özel dağılımlar hazırlanmamış, sınama süreçleri belirlenmemiş olmakla beraber çoğunluğunda, Pearson korelasyon katsayısı gibi sınama, doğruya yakın bir fikir verebilir. Geçerlik ve güvenirlilik katsayıları belirlenmesinde aksine bir belirtme yapılmadığı sürece, sözü geçen korelasyon katsayısı, ölçekten elde edilen verilerin yapısı uygun olması halinde Pearson korelasyon katsayısıdır.

1.2.2.1. Devamlılık Katsayısı (Coefficient of Stability). r_{cs}

Devamlı özellikler ile ilgili ölçmelerde aranan güvenirlilik, ölçümün devamlılığı açısındandır. Devamlılık katsayısı, ölçme aracının uygulanması, belli bir zaman aralığı sonunda aynı gruba aynı ölçme aracı uygulamasının tekrarlanması ile elde edilen iki seri halindeki sonuçlar arasındaki korelasyonun hesaplanması ile elde

edilir. Bu işlem "test-tekrar test" yöntemi olarak da bilinir. Devamlılık katsayısı, ölçme aracının bireyde kalıcı (devamlı) özellikleri ne ölçüde ölçtüğünü belirler.

Devamlılık katsayısının hesaplanmasında, iki uygulama arasında geçen sürenin ne kadar olması gerektiği sorunu ortaya çıkmaktadır. Bu soruya kesin bir cevap verilmemektedir. Bununla beraber, geçecek zaman, ölçeğin ölçtüğü özellik bakımından cevaplayıcıların önemli ölçüde değişmelerine yetecek kadar uzun olmamalıdır. Ayrıca, birinci uygulamada verilen cevapları hatırlamaya yetecek kadar kısa da olmamalıdır. Genellikle bu ara 3-6 hafta olarak belirlenebilir.

İki uygulamadan elde edilen puanlar arasındaki ilişki miktarının hesaplanması tek başına yeterli değildir. Bireylerin her iki ölçme aracı puanları arasında önemli bir fark olmamakla beraber her iki uygulamada cevaplandırılan, puan alınan maddeler farklı maddeler olabilir, aynı maddeden farklı puanlar almış olabilir. Bu nedenle, bireylerin her iki uygulamada puan aldıkları maddeler arasındaki tutarlılığa da bakmak gerekir. Test toplamı, faktör toplamı ve tüm maddeler için korelasyon hesaplanmalıdır.

Devamlılık katsayısı her ölçek için mutlaka hesaplanmak zorunda değildir. Ölçülen değişken devamlılık gösteren bir yapıdaysa devamlılık katsayısının bir anlamı vardır. Sürekli kaygı ölçüldüğünde devamlılık katsayısının belirlenmesi bir zorunluluk olmakla beraber, durumluk kaygının belirlendiği bir ölçekte devamlılık katsayısı anlamsızdır.

1.2.2.2. Eşdeğerlik Katsayısı (Coefficient of Equivalence). r_{ce}

Eşdeğerlik katsayısı, a) farklı uygulayıcıların aynı zamanda kullandıkları aynı ölçü araçları için **uygulamaların eşdeğerliğini** b) farklı, fakat benzer ölçme araçlarının (bu aynı değişkeni ölçen iki ayrı ölçek olabileceği gibi aynı ölçeğin iki ayrı formu da olabilir ve paralel test yöntemi olarak da bilinir) aynı gruba aynı zamanda uygulanmasında **araçların eşdeğerliğini**, c) aynı ölçme aracının farklı iki dildeki iki formunun aynı gruba aynı zamanda uygulanmasında **dilsel eşdeğerliği** belirtir. Bu uygulamalarda elde edilen puanlar arasındaki korelasyon katsayısı, eşdeğerlik katsayısıdır,

Paralel iki testin eşdeğerlik sınavının yapılabilmesi, her iki ölçme aracı, içindeki madde sayısı, niteliği, kullandığı ölçekleme tekniği (Likert, Thurstone, vb), faktör yapısı ve ölçtükleri davranış bakımından birbirine denk olmalıdır. Ölçek adaptasyon çalışmalarında bile bu eşdeğerlik koşullarına titizlikle uyulması gerekir. Örneğin orijinalinden bazı maddeleri hangi gerekçeyle olursa olsun ölçekten tamamen çıkararak ölçekteki madde sayısı azaltılırsa veya herhangi bir maddenin yerine bizce içerdiği değişken aynıdır diye yorumlanan yeni bir madde yazmak eşdeğerliği anlamsızlaştırır. Aynı şekilde orijinalinde "x", "q", "w" maddeleri "z" boyutu içinden yer alırken faktör analizi sonucunda bile olsa adapte ölçekte örneğin "q" maddesinin "g" boyutunda yer alması bu ölçeğin İngilizce-Türkçe formlarının eşdeğerlik sınavının önkoşulunu iptal eder. Bu iki form arasında ilişki bulunsa bile (ki büyük ihtimalle bulunacaktır) bu ilişki eşdeğerlik anlamında değildir.

Eşdeğer iki form, eğer aynı bireylere aynı zamanda uygulanmışsa eşdeğerlik katsayısı hesaplanır, ama belli bir zaman aralığı ile uygulanmışsa eşdeğerlik ve devamlılık birbirine karışır.

1.2.2.3. İçtutarlık Katsayısı (Coefficient of Internal Consistency), r_{ic}

Ölçme aracı bir gruba uygulandıktan sonra, sonuçların değerlendirilmesi aşamasında belli yöntemlere göre ikiye ayrılması ile bunlar arasında elde edilen katsayıya içtutarlık katsayısı denir. İçtutarlık katsayıları eğer ölçek tek bir boyuttan oluşursa ölçeğin bütününe göre yapılabileceği gibi ölçek faktörlerden oluşursa her faktör ve maddeleri kendi içinde bir bütün olarak kabul edilip aynı teknik yaklaşımla faktör düzeyinde de hesaplanabilir.

Tek bir ölçme aracı formu, tek bir birey' grubu ve tek bir uygulama gerektirdiğinden en sık kullanılan güvenilirlik saptama yöntemidir. Bu yöntemde karşılaşılan sorun, testin iki eşdeğer yarıya bölünmesidir. İki eşdeğer yarıya bölmenin en çok kullanılan yolu tek numaralı sorularla çift numaralı soruları ayrı puanlamaktır. Bu iki yarıdan elde edilen puanlar takımı, ayrı ayrı testlerden elde edilmiş gibi işlem görürler. Katsayı yüksek çıkmış, yani testin bir yansı diğer yarısı ile tutarlı ise, "bu testin içtutarlığı (güvenirliği) vardır" şeklinde değerlendirilecektir. Tek ve çift sorular şeklinde ikiye ayırmada göz önünde tutulması gereken nokta soruların numaralanmasında yanlılığa yol açabilecek herhangi bir düzenlemenin yapılmamış olması, soruların tek ve çift numaraları almasında random oluşmasıdır. Örneğin ölçek tek boyuttan oluşuyor ve tek numaralı sorular ölçülen değişken bakımından olumsuz (reverse) ifadeli maddeler çift numaralı sorular ise olumlu ifadeli maddeler ise bu ikiye ayırma tekniği doğru sonuçlar vermeyebilir.

İçtutarlık katsayısını hesaplamada, bu temel yaklaşımı kullanmakla beraber farklı istatistik değer ve formülleri kullanan çok sayıda teknik vardır. Veri türleri ve koşullara uygun oldukları sürece, güvenilirlik ile ilgili sınamaları pekiştirmek amacıyla bu tekniklerin birden fazlası kullanılabilir.

İçtutarlık katsayılarından Spearman-Brown, Stanley, Rulon, Flanagan, Mossier, Horst gibi bazı teknikleri yan-test (split-half) tutarlık katsayıları olarak adlandırılır. İtem-total korelasyon katsayısı, item-remainder korelasyon katsayısı teknikleri ise parça-bütün arasındaki ilişki yoluyla tutarlığı belirler, her madde için ayrı ayrı hesaplanır. Kuder-Richarson 20 ve 21 ile Cronbach ∞ tekniği ise parçalar arası ortak ilişkiyi dikkate alarak bütün için tek bir tutarlık katsayısı hesaplar.

1.2.2.3.1. Spearman-Brown

Bir gruba uygulanmış testin iki eşdeğer yarıya bölünmesi ve bireylerin iki eşdeğer yarıdan aldıkları puanlar arasındaki Pearson korelasyon katsayısından hareketle Spearman-Brown formülünden de yararlanarak testin bütününe güvenilirliği kestirilir.

Spearman-Brown formülü,

r_{ic} : testin güvenilirlik katsayısı

r_{xy} : iki yarı arasındaki korelasyon katsayısı
olarak ifade edilirse

$$r_{ic} = \frac{2r_{xy}}{1 + r_{xy}} \quad [2.5]$$

olur.

Tablo 3'deki verilere göre yan testin güvenilirliği 0.468 olarak bulunmuştur. Bu bulgu, testin bütününe güvenilirliğini bulmak için Spearman/Brown formülüne [2.5] yerleştirildiğinde

$$r_{ic} = \frac{2 * (0.468)}{1 + 0.468} = 0.638$$

bulunur.

TABLO 3. SPEARMAN-BROWN TEKNİĞİ İLE İÇTUTARLIK KATSAYISI HESAPLANIŞI

Öğrenci	Tek	Çift	X^2	Y^2	XY
	Numaralı	Numaralı			
	Sorular	Sorular			
	Toplamı	Toplamı			
	X	Y			
A	20	20	400	400	400
B	20	20	400	400	400
C	18	17	324	289	306
D	24	9	576	81	216
E	15	15	225	225	225
F	14	15	196	225	210
G	14	14	196	196	196
H	17	9	289	81	153
I	8	8	64	64	64
İ	9	7	81	49	63
10	159	134	2751	2010	2233
n	S_x	S_y	S_x^2	S_y^2	S_{xy}

$$r_{xy} = \frac{(10 * 2233) - (159 * 134)}{\sqrt{[10 * 2751 - (159)^2] * [(10 * 2010) - (134)^2]}} = 0.468$$

1.2.2.3.2. Stanley

Bir testin iki yarısındaki puanlardan test güvenilirliğini kestirmede kullanılan Stanley formülü;

$\Sigma\ddot{u}x-y$: üst çeyreklikte yer alanların x-y toplamı

$\Sigma ax-y$: alt çeyreklikte yer alanların x-y toplamı

$\Sigma\ddot{u}x+y$: üst çeyreklikte yer alanların x+y toplamı

$\Sigma ax+y$: alt çeyreklikte yer alanların x+y toplamı

$$r_{ic} = 1 - \frac{[(\Sigma\ddot{u}x - y) - (\Sigma ax - y)]^2}{[(\Sigma\ddot{u}x + y) - (\Sigma ax + y)]^2} \quad [2.6]$$

Ölçeğin güvenilirliğini kestirmede uygulanan toplam test sayısının $nq = (n \cdot 0,27)$ bulunur. Ölçeğin her iki yarısındaki puanların toplamına (X+Y) göre üst çeyreklikte yer alanların (\ddot{u}) puanlarının toplamı ($\Sigma\ddot{u}x+y$) ve alt çeyreklikte yer alanların (a) puanlarının toplamı ($\Sigma ax+y$) bulunur. İşlem benzer olarak ölçeğin her iki yarısı arasındaki farklara (X-Y) göre üst çeyreklikte yer alanların puanlarının toplamı ($\Sigma\ddot{u}x-y$) ve alt çeyreklikte yer alanların puanlarının toplamı ($\Sigma ax-y$) bulunur. Elde edilen bu dört değer Stanley formülündeki [2.6] yerine konularak bu yöntemle bu yönteme göre içtutarlık belirlenmiş olur.

TABLO 4. STANLEY TEKNİĞİ İLE İÇTUTARLIK KATSAYISI HESAPLANIŞI

Öğrenci	TEK X	ÇİFT Y	(3) X+Y	(4) X-Y
A	20	20	40 \ddot{u}	0a
B	20	20	40 \ddot{u}	0a
C	18	17	35 \ddot{u}	1
D	24	9	33	15 \ddot{u}
E	15	15	30	0
F	14	15	29	-1a
G	14	14	28	0
H	17	9	26a	8 \ddot{u}
I	8	8	16a	0
İ	9	7	16a	2 \ddot{u}
10			$\Sigma\ddot{u}x+y=115$	$\Sigma\ddot{u}x-y=24$
n			$\Sigma ax+y= 58$	$\Sigma ax-y=- 1$

Tablo 3'deki birimlere ait aynı verileri kullanarak hazırlanan Tablo 4 için bu yöntem uygulandığında $nq = (10 \cdot 0,27) = 3$ kişi üst çeyreği 3 kişi alt çeyreği oluşturmaktadır. Tablonun altında hesaplanmış şekilde verilen değerler formüle [2.6] uygulandığında

$$r_{ic} = 1 - \frac{(24 - (-1))^2}{(115 - 24)^2} = 1 - \frac{25^2}{91^2} = 1 - \frac{625}{8281} = 0,924$$

olarak elde edilir.

1.2.2.3.3. Rulon

Bir testin iki yarısındaki puanlardan test güvenilirliğini kestirmede kullanılan Rulon formülü,

TABLO 5. RULON, FLANAGAN, MOSIER VE HORST TEKNİKLERİ İLE İÇTUTARLIK HESAPLANIŞI

Öğrenci	X	Y	X+Y	X-Y
A	20	20	40	0
B	20	20	40	0
C	18	17	35	1
D	24	9	33	15
E	15	15	30	0
F	14	15	29	-1
G	14	14	28	0
H	17	9	26	8
I	8	8	16	0
i	9	7	16	2
$\sum X$	159	134	293	25
$\sum X^2$	2751	2010	9227	295
μ	15.90	13.40	29.30	2.50
α	4.98	4.88	8.45	5.08
α^2	24.77	23.82	71.34	25.83
n=10	$r_{xy} = 0.468$	$r_{xt} = 0.860$	$r_{yt} = 0.854$	

Tabloda X (çift numaralı soruları), Y (tek numaralı soruları), X+Y (test toplam puanını), X-Y (her iki yarı arasındaki farkı) ifade etmektedir

Daha önceki tekniklerde de kullanılan verileri içeren Tablo 5'deki sonuçlara göre

$$r_{ic} = 1 - \frac{25.83}{71.34} = 1 - 0.36 = 0.638$$

olur.

1.2.2.3.4. Flanagan

Bir testin iki yarısındaki puanlardan test güvenilirliğini kestirmede kullanılan Flanagan formülü, σ

σ^2_x : Tek numaralı soruların varyansı

σ^2_y : Çift numaralı soruların varyansı

σ^2_t : toplam puanların varyansı

olarak ifade edilirse

$$r_{ic} = 2 * \left(1 - \frac{\sigma_x^2 + \sigma_y^2}{\sigma_t^2} \right) \quad [2.8]$$

Tablo 5'deki sonuçlara göre

$$r_{ic} = 2 * \left(1 - \frac{24.77+23.82}{71.34} \right) = 2 * (1-0.68) = 0.640$$

olur.

1.2.2.3.5.Mosier

Mosier formülü,

r_{xt} : tek numaralı sorularla toplam arasındaki korelasyon katsayısı

σ_x^2 : Tek numaralı soruların varyansı

σ_t^2 : Toplam puanların varyansı

olarak ifade edilirse

$$r_{ic} = \frac{r_{xt}\sigma_t - \sigma_x}{\sqrt{[\sigma_t^2 + \sigma_x^2 - (2r_{xt}\sigma_x\sigma_t)]}} \quad [2.9]$$

Tablo 5'deki sonuçlara göre

$$r_{ic} = \frac{0.860 * 8.447 - 4.977}{\sqrt{[71.344 + 24.767 - (2 * 0.860 * 4.977 * 8.447)]}} = \frac{2.287}{4.879} = 0.468$$

olur.

Elde edilen değer Spearman-Brown formülüne [2.5] yerleştirilerek testin bütünüünün güvenilirliği bulunur.

1.2.2.3.6.Horst

Şimdiye kadar bahsedilen içtutarlık katsayıları hesaplanan iki varırım eşit sayıda maddeden oluştuğu varsayımına dayanır. Farklı sayıda maddelerden oluşan iki parça arasındaki tutarlığı hesaplamak için buna yönelik düzeltmeleri içeren Horst tekniği kullanılır. Güvenirlik hesabı yapılan ölçeğin iki yansının eşit olmadığı durumlarda kullanılabilmesi bu tekniğin ayırıcı özelliğidir. Bu teknik doğrudan iki yarı arasındaki içtutarlığı hesaplamada kullanılabilceği gibi esas olarak, testin bir alt bölümünün kendi hariç diğer alt bölümlerin toplamına göre tutarlı olup olmadığını sınamada kullanılır. Eğer iki parçadaki madde sayıları eşit ise diğer tekniklerle varılan sonucu verir, ancak madde sayıları arasındaki denge bozuldukça diğer içtutarlık katsayılarından ve korelasyondan farklılık gösterebilir.

Horst formülü,

- r_{xy} : iki yarı arasındaki korelasyon katsayısı
 p : Birinci yarıdaki soru sayısının toplam soru sayısına oranı
 q : $1-p$

Olarak ifade edilirse

$$r_{ic} = \frac{r_{xy} \{ \sqrt{[(r_{xy}^2 + 4pq(1-r_{xy}^2))] - r_{xy}} \}}{2pq(1-r_{xy}^2)} \quad [2.10]$$

Tablo 5'deki sonuçlara göre

$$r_{ic} = \frac{0.468 [(\sqrt{0.468^2 + 4 (0.50) (0.50) (1-0.468^2)}) - 0.468]}{2(0.50) (0.50) (1-0.468^2)}$$

$$r_{ic} = \frac{0.468 [(\sqrt{1}) - 0.468]}{0.390} = \frac{0.468 * 0.532}{0.390} = \frac{0.249}{0.390} = 0.638$$

olur.

1.2.2.3.7. Cronbach α

Bu tekniğin özelliği Kuder-Richardson tekniğinkine benzer. Testin tüm alt bölümlerinin birbirlerine (hepsi dahil toplama) göre, ya da bir alt bölümün tüm sorularının birbirlerine (hepsi dahil toplam) göre, tutarlı olup olmadığını sınamada kullanılır. Sonuçlar tüm bölümlerin (ya da o bölümdeki tüm soruların) birbiriyle ilgili olduğu şeklinde yorumlanır. Katsayının yüksekliği içtutarlığın yüksekliğini gösterir.

Bir bölümdeki soruların o bölümle ilgisi söz konusu olduğunda Croanbach α formülü,

m : bölümdeki soru sayısı
 σ_j^2 : j. sorunun varyansı
 σ_t^2 : soruların toplamının varyansı
 $\Sigma\sigma_j^2$: soruların varyanslarının toplamı (1. sorunun varyansı +.....+ j. sorunun varyansı)

$$\alpha = \frac{m}{m-1} * (1 - \frac{\Sigma\sigma_j^2}{\sigma_t^2}) \quad [2.11]$$

olur.

Aynı formüldeki semboller, alt bölümlerin testin bütünüyle ilgisi söz konusu olduğunda,

m : testteki altbölüm sayısı

σ_j^2 : j. alt bölümün toplam puanlarının varyansı

σ_t^2 : bölüm toplamalarının toplamının varyansı

$\Sigma\sigma_j^2$: bölüm toplam varyanslarının toplamı

anlamında kullanılır.

TABLO 6. CRONBACH, ITEM-TOTAL KORELASYON, ITEM-REMAINDER KORELASYON HESAPLANIŞI

Öğrenci	I ₁ 2.soru	I ₂ 4.soru	I ₃ 6. soru	T toplam	T _{I₁} 2.s * tp
A	3	3	5	11	33
B	4	3	5	12	48
C	3	3	4	10	30
D	2	2	1	5	10
E	4	2	5	11	44
F	3	3	3	9	27
G	4	2	2	8	32
H	2	2	1	5	10
I	1	1	2	4	4
İ	1	1	2	4	4
Σx	27	22	30	79	242
Σx^2	85	54	104	713	
μ	2.70	2.20	3.00	7.90	
σ	1.160	0.789	1.633	3.143	
σ^2	1.344	0.622	2.667	9.878	

Tablo 6'deki sonuçlara göre, 2., 4. ve 6. sorulardan oluşan bir alt ölçek için

$$\alpha = \frac{3}{3-1} * (1 - \frac{(1.344+0.622+2.667)}{9.878}) = 1.5 * (1-0.469) = 0.796 \text{ olur.}$$

1.2.2.3.8. Kuder-Richardson 20-21

Testteki herbir maddeye doğru cevap veren birey yüzdesi hesaplanarak içtutarlık Kuder-Richardson 20 ve 21 no'lu formülleri ile tahmin edilebilir. Bu yolla elde edilen katsayı da iç-tutarlı bir ölçüdür. "

Kuder-Richardson (KR) 20 no'lu formülü şöyledir:

m : testteki madde sayısı

P : bir maddeyi doğru cevaplayanların oranı

q : bir maddeyi doğru cevaplamayanların oranı = (1-p)

$\sum pq$ her madde için hesaplanan (p x q)'ların toplamı

$\alpha^2 t$: test toplam puanlarının varyansı

olarak ifade edilirse

$$r_{ic} = \left(\frac{m}{m-1} \right) * \left(\frac{\sigma_t^2 - \sum pq}{\sigma_t^2} \right) \quad [2.12]$$

olur.

Bu formül [2.12], doğru cevaplandırılanlara puan verilmesi, yanlış ve boş cevaplara hiç puan verilmediği testlerde uygulanabilir. Cevapların doğruluğu-yanlışlığı mantığına oturduğu için başarı testlerinde kullanılabilmesine karşılık tutum testlerinde kullanılmamalıdır. Eğer testteki maddeler farklı ağırlıklarla puanlanmışsa bu formül uygulanamaz. Bir testteki maddelerin güçlük dereceleri birbirinden önemli ölçüde farklı değilse o testin güvenilirliği için KR-21 no'lu formülü kullanılabilir. Bu formül:

p : madde doğru cevap oranlarının (p) aritmetik ortalaması

q : q oranların aritmetik ortalaması

$$r_{ic} = \left(\frac{m}{m-1} \right) \left(\frac{\sigma_t^2 - m\bar{p}\bar{q}}{\sigma_t^2} \right) \quad [2.13]$$

olur. Toplam puanların ortalaması X_t gösterilirse $p = X_t / m$ olduğundan ve $q = (m-X_t) / m$ olduğundan formül [2.13]

$$r_{ic} = \frac{m\sigma_t^2 - \bar{X}_t (m - \bar{X}_t)}{(m-1)\sigma_t^2} \quad [2.14]$$

şeklini alır. KR-21 içtutarlığı kestirmede (özellikle de KR-20'ye göre daha) düşük güçtedir.

1.2.2.3.9. İtem-Total Korelasyon Katsayısı r_{it}

Bu yöntemde her bir sorudan elde edilen puan ile toplam puan arasındaki ilişkiye bakılır Pearson Çarpım Moment korelasyon katsayısı yoluyla her bir madde için elde edilen sonuçların manidarlığına bakılır.

$$r_{ii} = \frac{n \sum T_i - (\sum T \sum i)}{\sqrt{n \sum T^2 - (\sum T)^2} \sqrt{n \sum i^2 - (\sum i)^2}} \quad [2.15]$$

- $\sum T_i$: her bireyin toplamdaki ve itemdeki puanların çarpımlarının tüm bireyler için toplamı
- $\sum T$: bireylerin toplam puanlarının toplamı
- $\sum i$: bireylerin itemdeki(maddedeki) puanlarının toplamı
- $\sum T^2$: bireylerin toplam puanlarının karelerinin toplamı
- $\sum i^2$: bireylerin itemdeki (maddedeki) puanlarının karelerinin toplamı
- n : birey sayısı

Tablo 6'daki verilerden 2 nolu soru için r_{ii} hesaplamak için veriler formül [2.15]'e

$$r_{ii} = \frac{10 \times 242 - (79 \times 27)}{\sqrt{10 \times 713 - 79^2} \sqrt{10 \times 85 - 27^2}} = 0.875 \quad \text{uygulandığında;}$$

olarak bulunur. Hesaplama formülü [2.15] doğrudan doğruya Pearson formülü olduğundan bu istatistiğe göre manidarlığı sınıandığında hesaplanan $r=0.875$, $sd(\text{serbestlik derecesi}) = (n-2) = 8$ için tablo değeri olan 0.765'den büyük olduğu için .01 düzeyinde manidardır. Bu madde, bu ölçekle ilişkilidir. Bu maddenin ölçtüğü değişken ile ölçeğin tümünün ölçtüğü değişken birbiriyle ilişkilidir. O halde bu madde bakımından içtutarlık vardır. Manidar çıkmayan madde diğer analizler sonuçlarına bakılmaksızın ölçekten atılır. Bazen r_{ii} manidar fakat negatif bir değer olarak çıkabilir. Bu durumda yapılması gereken ayırt edicilik için t-testi açıklanırken belirtildiği gibi maddenin değerlendirme yönünü değiştirmek ve tabii ki geçerlik güvenilirlik işlemlerinin yenilemek. Maddenin ölçtüğü değişken ile ait olduğu ölçeğin ölçtüğü değişken "ilişkilidir" ifadesinin kullanılmasına dikkat etmek gerekir, "Aynıdır" ifadesi kullanılamaz, aynı olduğunu ileri sürebilmek için öncelikle iki ölçüm (madde-ölçek) puanları arasında tam bir ilişki ($r=+1$) olmalı, ayrıca fark olmadığı bağlantılı grup t-testi ile sınıanarak fark bulunmamalıdır.

1.2.2.3.10. Item Remainder Korelasyon Katsayısı r_{ri}

Bu tekniğin özü (j. soru) ile (j. soru hariç ölçeğin tümünden alınan puan) arasındaki ilişkiyi, bir maddenin kendi hariç ait olduğu ölçeğin toplamla ilişkisini bulmaktır.

Item-remainder formülü,

r_{ti} : j. soru ile toplam puan arasındaki korelasyon katsayısı

σ_j^2 : j. sorunun varyansı

σ_t^2 : Toplam puanların varyansı

olarak ifade edilirse

$$r_{ti} = \frac{r_{ti} \sigma_t - \sigma_i}{\sqrt{\sigma_t^2 + \sigma_i^2 - 2r_{ti} \sigma_t \sigma_i}} \quad [2.16]$$

Tablo 6'daki verilerden 2 nolu soru için r_{ti} hesaplamak için veriler formüle [2.16] uygulandığında

$$r_{ti} = \frac{0.875 \times 3.143 - 1.160}{\sqrt{9.878 + 1.344 - 2 \times 0.875 \times 3.143 \times 1.160}} = 0.723$$

Hesaplanan r_{ti} tablo değeri 0.632'den büyük olduğu için .05 düzeyinde manidardır. Bu madde ölçeği oluşturan diğer 2 madde ile ilişkilidir, ölçtükleri değişkenler ilişkilidir.