

# Krosstöflur

Fyrirlestur í Tölfræði II (SÁL203G)

# 2×2 töflur

Myndin sýnir fjölda þeirra sem detta í það eftir kyni.

Við lítum á kynferði sem frumbreytu og látum hana því skilgreina dálkana.

Þetta eru tvær tvíkostabreytur, taflan hefur því fjögur hólf (*cells*). Svona töflur eru kallaðar 2×2 töflur.

Heildartölur fyrir kynferði og það að detta í það kallast jaðardreifingar. Ef við skoðum tölurnar fyrir hvort kyn fyrir sig, eru það skilyrtar dreifingar (*conditional distribution*).

Detta í það	Kynferði		Samtals
	Karlar	Konur	
Já	1.630	1.684	3.314
Nei	5.550	8.232	13.782
Samtals	7.180	9.916	17.096

Að detta í það er fylgibreyta og skilgreinir línur.

Kynferði er frumbreyta og skilgreinir dálka.

Skilyrt dreifing fyrir karla.

Þetta er eitt af hólfum töflunnar.

Jaðardreifing

# Sameiginleg og skilyrt dreifing

Tölur eru oft birt sem prósentur. Við getum tekið prósentur af heildarfjölda í töflunni og fengið sameiginlega (*joint*) dreifingu breytanna tveggja

Við getum einnig reiknað prósentur innan hvers dálks og fengið þannig skilyrtar dreifingar.

Að síðustu má reikna prósentur innan hvernar línu en það er oft órökrétt þar sem þá fæst skilyrt dreifing miðað við *fylgibreytuna*.

Gott er að gefa upp fjöldann sem prósentur miðast við.

Detta í það	Kynferði		Samtals
	Karlar	Konur	
Já	9,5%	9,9%	19,4%
Nei	32,5%	48,2%	80,6%
Samtals	42,0%	58,0%	100,0%

Athugasemd. Taflan sýnir sameiginlega dreifingu (*joint distribution*).  $N= 17.096$ .

Detta í það	Kynferði		Samtals
	Karlar	Konur	
Já	22,7	17,0	19,4
Nei	77,3	83,0	80,6
Samtals	100,0	100,0	100,0
<i>n</i>	7.180	9.916	17.096

Skilyrtar dreifingar er oft sýnd í prósentum innan hvers dálks.

Jaðardreifingu má einnig skoða í prósentum.

# Áhrif tónlistar á vínsölu

Töflurnar sýna áhrif tónlistar á sölu léttvíns. Í írskum stórmarkaði var ýmist spiluð frönsk, ítölsk eða engin tónlist og athugaður fjöldi seldra flaskna.

Efri taflan sýnir fjöldataölur og bera með sér að frönsk og önnur vín eru vinsælust. Frönsk tónlist dregur úr sölu ítalskra og eykur sölu franskra vína; ítölsk tónlist eykur sölu ítalskra vína.

Skilyrtu dreifingarnar á neðri myndinni auðvelda túlkunina.

Vín	Tónlist			Heild
	Frönsk	Ítölsk	Engin	
Frönsk	39	30	30	99
Ítölsk	1	19	11	31
Önnur	35	35	43	113
Heild	75	84	84	243

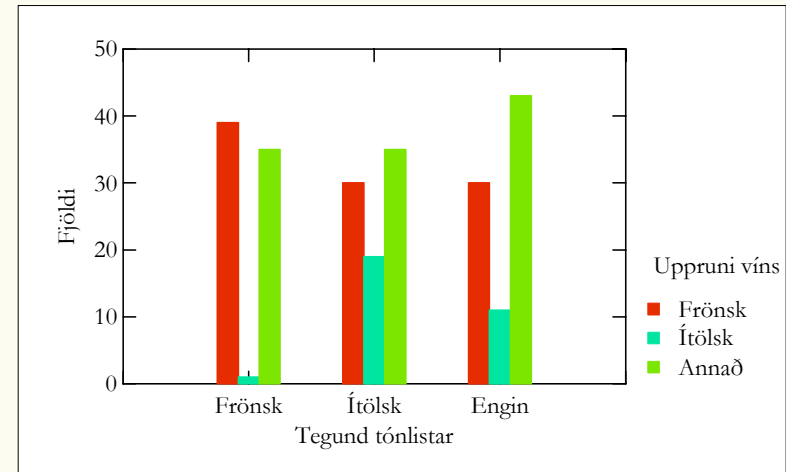
Vín	Tónlist			%
	Frönsk	Ítölsk	Engin	
Frönsk	52,0	35,7	35,7	40,7
Ítölsk	1,3	22,6	13,1	12,8
Önnur	46,7	41,7	51,2	46,5
Heild	100,0	100,0	100,0	100,0
<i>n</i>	75	84	81	243

# Hvernig birtast tengsl í krosstöflu?

Ef engin tengsl eru milli tónlistar og vínsölu, ættu allar skilyrtar dreifingar að vera eins.

Myndin sýnir skilyrtar dreifingar fyrir ólíkan uppruna tónlistar. Við horfum á *hlutföllin* þar sem fjöldinn getur verið ólíkur eftir því hver tónlistin er.

Mynstur vínsölu er greinilega mjög ólíkt eftir tónlist. Frönsk tónlist sker sig úr vegna hlutfallslega lítillar sölu á ítölskum vínum. Ítölsk tónlist er samfara auknu hlutfalli ítalskra vína, á kostnað annarra vína.



Yfirleitt er best að nota fjöldatölur í súlurítum og láta frumbreytuna skilgreina klasana eins og hér er gert. Fjöldatölur gefa allar helstu upplýsingar en sjónskynjunin er næm á hlutföllin. Það má birta hlutföllin beint—eins og Moore og McCabe—en það er síðra.

# Núlltilgátan

Það er engin formleg núlltilgáta með táknum fyrir kíkvaðratpróf, heldur kveður hún einfaldlega á um að engin tengsl séu í töflunni.

Í vindaðminu birtist tengslaleysi sem nákvæmlega eins skilyrtar dreifingar, þ.e. 40,7% sölunnar séu frönsk, 12,8% ítölsk og 46,5% önnur vín burt séð frá því hvaða tónlist kann að vera spiluð.

Þannig fæst væntitíðni, fjöldinn í hverju hólfi ef hlutföllin væru þau sömu í öllum dálkum.

Vín	Tónlist			Heild
	Frönsk	Ítölsk	Engin	
Frönsk	39	30	30	99
Ítölsk	1	19	11	31
Önnur	35	35	43	113
Heild	75	84	84	243

Rauntíðni

Vín	Tónlist			Heild
	Frönsk	Ítölsk	Engin	
Frönsk	30,6	34,2	34,2	99,0
Ítölsk	9,6	10,7	10,7	31,0
Önnur	34,9	39,1	39,1	113,0
Heild	75,0	84,0	84,0	243,0

Væntitíðni

# Kíkvaðratpróf

Ef tíðnin er lík væntitíðninni er frávikið frá núlltilgátunni lítið. Við viljum vita hvort frávikið sé það mikið að við treystum okkur til að hafna núlltilgátunni.

Kíkvaðratprófið ber saman raun- og væntitíðni og er mælikvarði á mismun þeirra. Ef niðurstaða prófsins er há, er mikið frávik frá núlltilgátunni. Ef frávikið er það mikið að líkurnar á þetta miklum eða meiri frávikum eru  $\alpha$  eða lægri, getum við hafnað núlltilgátunni.

$$\chi^2(df, N = fjöldi) = \sum \frac{(rauntíðni - væntitíðni)^2}{væntitíðni}$$

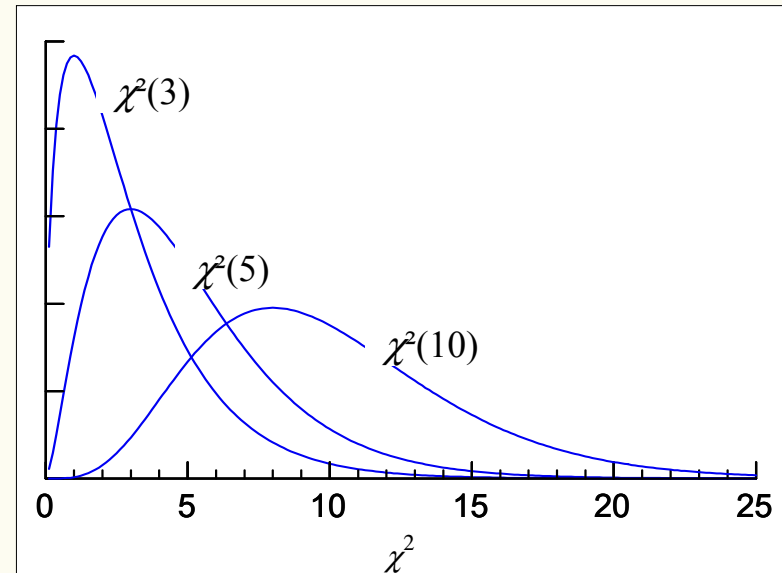
Þetta próf er sambærilegt z- eða t-prófi þar sem niðurstaðan gefur frávikið frá  $H_0$ . Ólíkt þeim er það alltaf stefnulaust. Niðurstöðunni flettum við upp í töflu eða berum saman við vendigildi. Ef við fáum niðurstöðu sem er jafnhá eða hærri en vendigildið sem við finnum í töflu F, höfnum við núlltilgátunni.

# Kíkvaðratdreifing

Kíkvaðrat er fjölskylda dreifinga eins og  $t$ -dreifing. Lögun dreifingarinnar fer eftir frígráðunum.

Niðurstaðan verður aldrei lægri en 0, skekkjan er ætíð jákvæð en minnkar með auknum frígráðum. Meðaltal dreifingarinnar er jöfn frígráðunum.

Myndin sýnir dæmi um þrjár dreifingar. Það er greinilegt að með auknum frígráðum verður dreifingin samhverfari.





# Forsendur kíkvaðratprófs

Kíkvaðratpróf gefur nálgun að úrtakadreifingunni sem batnar eftir því sem fjöldinn er meiri. Það er núlltilgátan sem skiptir máli, þ.e. væntitíðnin.

Algengt viðmið er að ekkert hólf hafi væntitíðni undir 1,0 og í mesta lagi 20% þeirra séu undir 5,0. M&M miða við að *meðaltal* væntitíðni sé 5 eða herra.

Mælingar þurfa að vera óháðar en það er uppfyllt ef hvert stak kemur aðeins einu sinni fyrir í töflunni.

## Forsendur kíkvaðratprófs

- Óháðar mælingar
- Mælingar byggjast á tíðni
- Væntitíðni 1,0 eða hærri í öllum hólfum
- Væntitíðni undir 5 í mest 20% hólfanna

# Þegar forsendur brestur

Ef mælingar byggjast ekki á tíðni, hef ég augljóslega beygt einhvers staðar út af (NB! Ekki á gatnamótum þó, er utan vega). Væntanlega er einhver önnur úrvinnsla sem hentar. Ef t.d. við erum með meðaltalatöflu, gæti  $t$ -próf eða dreifigreining hentað.

Ef mælingar eru háðar, hefur okkur orðið á í rannsóknarsniðinu. Við því er fátt að gera annað en að byrja upp á nýtt.

Ef væntitíðnin er of lág en taflan er  $2 \times 2$ , get ég notað Fishers Exact Test. Það gefur mér nákvæm líkindi undir núlltilgátunni án þess að gera kröfur til væntitíðninnar.

Að öðrum kosti get ég hugleitt að fella saman flokka hjá annarri hvorri breytunni eða báðum. Ég get jafnvel hugleitt að fella brott einhverja flokka.

Slíkar *breytingar* á töflunni þarf þó að íhuga vel og gæta þess að þær séu fullkomlega réttlætanlegar.

# Samanburður við z-próf á hlutföll

Kíkvaðratpróf er náskyld z-prófi fyrir mismun tveggja óháðra hlutfalla.

Þar notuðum við hlutfallið 0,194 sem mat á hlutfallinu undir núlltilgátunni. Við athugðum síðan hvort að munur hlutfallanna 0,227 og 0,170 væri það mikill að ólíklegt væri að hlutföllin endurspegluðu sameiginlega hlutfallið 0,194.

Kíkvaðratpróf notar hlutfallið 0,194 til að gera væntitíðni og ber rauntíðni saman við væntitíðnirnar.

Detta í það	Kynferði		Samtals	%
	Karlar	Konur		
Já	1.630	1.684	3.314	19,4
Nei	5.550	8.232	13.782	80,6
Samtals	7.180	9.916	17.096	100,0

Það er stærðfræðileg tengsl milli kíkvaðratprófs og z-próf fyrir hlutföll.

Kíkvaðratpróf gefur  $\chi^2(1, N= 17.096)= 87,172$ . Kvaðratrótin af því er 9,337 sem er jafnt niðurstöðu z-prófsins.

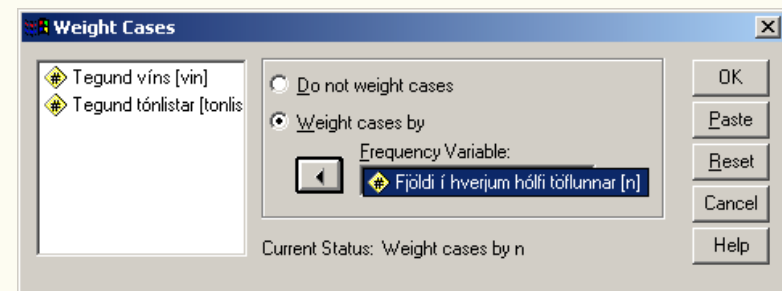
# Unnið með vigtuð gögn í SPSS

Ég les inn skjalið **Tónlist og vínsala.sav**. Í því er ein færsla fyrir hvert hólf í töflunni og breytan *N* gefur til kynna fjöldann í hólfinu.

Ég þarf því að nota hverja færslu oft, jafn oft og fjöldinn er í hverju hólfu.

Ég fer í **Data/Weight Cases...** og vel þar breytuna *N* eins og myndin sýnir. Smelli síðan á **OK**. Við það verður *N* eins konar margfaldari, þ.e. notar hverja færslu jafnoft og gildi breytunnar segir til um.

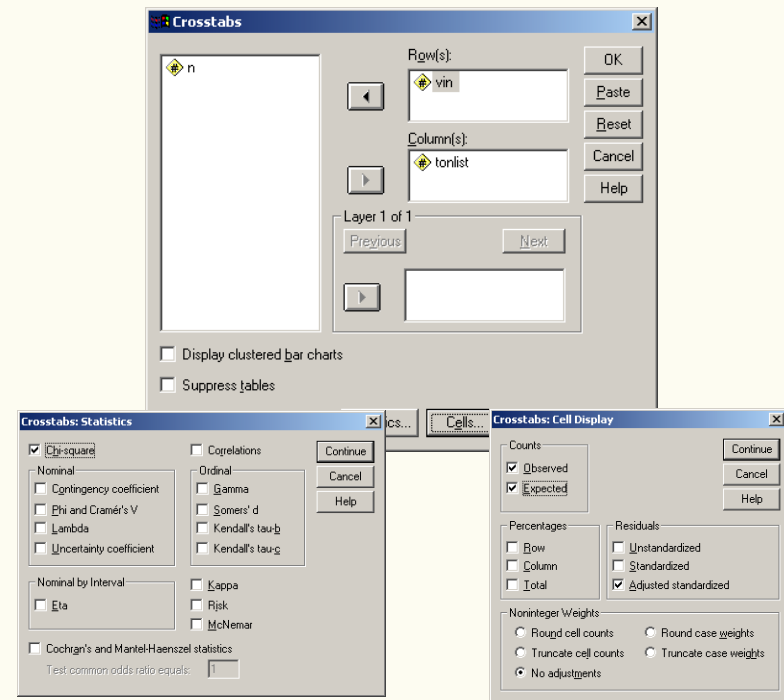
Tonlist	Vin	N
Frönsk	Frönsk	39
Ítölsk	Frönsk	30
Engin	Frönsk	30
Frönsk	Ítölsk	1
Ítölsk	Ítölsk	19
Engin	Ítölsk	11
Frönsk	Annað	35
Ítölsk	Annað	35
Engin	Annað	43



# Kíkvaðratpróf í SPSS

Ég fer í Analyze /Descriptive Statistics /Crosstabs... færi Vin í efri og Tonlist í neðri textareitinn. Ég smelli síðan á takkann Statistics og haka þar við Chi square. Undir Cells... haka ég við Observed, Expected og Adjusted standardized.

Að þessu loknu smelli ég á OK. Við það birtast niðurstöðurnar í niðurstöðuglugganum.



# Niðurstöður í SPSS

Í niðurstöðuglugganum skoðum við fyrst fjöldatölurnar í töflunni og göngum úr skugga um að þær séu réttar eða a.m.k. sannfærandi.

Síðan er niðurstaða kíкваðratprófsins skoðuð. Við fáum  $\chi^2(4, N= 243)= 18,3, p= 0,001$ . Við getum því hafnað núlltilgátunni og dregið þá ályktun að vínsala tengist þeirri tónlist sem er spiluð. Ef núlltilgátan væri rétt, væru minna en 0,1% líkur á svona miklu eða meira frávik frá  $H_0$  í úrtaki.

vin Tegund vins \* tonlist Tegund tónlistar Crosstabulation

		tonlist Tegund tónlistar			Total	
		Frönsk	Ítölsk	Engin		
vin Tegund vins	Frönsk	Count	39	30	30	99
		Expected Count	30,6	34,2	34,2	99,0
		Adjusted Residual	2,4	-1,2	-1,2	
	Ítölsk	Count	1	19	11	31
		Expected Count	9,6	10,7	10,7	31,0
		Adjusted Residual	-3,6	3,3	,1	
	Annað	Count	35	35	43	113
		Expected Count	34,9	39,1	39,1	113,0
		Adjusted Residual	,0	-1,1	1,1	
Total	Count	75	84	84	243	
	Expected Count	75,0	84,0	84,0	243,0	

Chi-Square Tests

	Value	df	Asymp. Sig. (2-sided)
Pearson Chi-Square	18,279 <sup>a</sup>	4	,001
Likelihood Ratio	21,875	4	,000
Linear-by-Linear Association	1,962	1	,161
N of Valid Cases	243		

a. 0 cells (.0%) have expected count less than 5. The minimum expected count is 9.57.

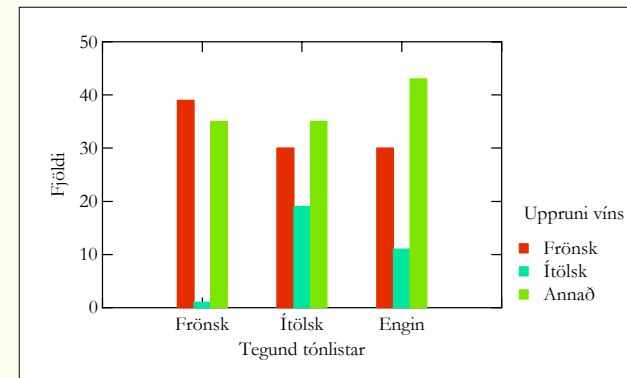
*Lestu alltaf þessa neðanmálsgrein!*

# Efnisleg túlkun sambandsins

Ef ályktun okkar er rétt, þá eru tengsl milli tónlistar og vínsölu. Við vitum hins vegar ekki *hver* tengslin eru.

Við höfum þrjár leiðir til að túlka niðurstöðuna. Við getum (a) skoðað töfluna eða súlurit og séð þannig hvers eðlis tengslin eru, (b) borið saman raun- og væntitíðni og (c) túlkað leiðréttu leif.

Leiðrétt leif er á kvarða staðalvillu og því er niðurstaða sem er hærri en 1,96 eða lægri en  $-1,96$  óvenjuleg.



vin Tegund víns \* tonlist Tegund tónlistar Crosstabulation

		tonlist Tegund tónlistar			Total	
		Frönsk	Ítölsk	Önnur		
vin Tegund víns	Frönsk	Count	39	30	30	99
		Expected Count	30,6	34,2	34,2	99,0
		Adjusted Residual	2,4	-1,2	-1,2	
	Ítölsk	Count	1	19	11	31
		Expected Count	9,6	10,7	10,7	31,0
		Adjusted Residual	-3,6	3,3	,1	
	Annað	Count	35	35	43	113
		Expected Count	34,9	39,1	39,1	113,0
		Adjusted Residual	,0	-1,1	1,1	
Total		Count	75	84	84	243
		Expected Count	75,0	84,0	84,0	243,0

# Kíkvaðratpróf í CrunchIt

Fyrst set ég töfluna upp í CrunchIt eins og sést á myndinni eða sæki **Tónlist og vínsala (enskt).csv**. Textinn þarf að vera á ensku.

Síðan fer ég í **Statistics/Tables/Contingency**. Í glugganum sem birtist vel ég **with counts**, vel dálka töflunnar og gef upp hvaða dálkur hefur fjöldann.

Niðurstaðan er ekki eins ítarleg og í SPSS, t.d. er leifin ekki birt og hlutföll eru gefin upp miðað við hlutfall af heildarfjölda.

The screenshot shows the CrunchIt 2.0 interface. On the left is a navigation menu with 'Data', 'Statistics', 'Summary Statistics', 'Tables', 'Frequency', 'Contingency', 'Z Test', 'Prop Test', 'T Test', 'Variation', 'Regression', and 'ANOVA'. The main window displays a data table with 8 rows and 4 columns. A 'Contingency Table' dialog box is open, showing 'vine' and 'music' as the first and second variables, and 'n' as the counts variable. The 'With Counts' radio button is selected. Below the dialog, a 'Contingency Table with Counts' table is shown, along with a 'Chi-Square' table.

Row	Col 1	Col 2	Col 3	Col 4
#	vine	music	n	
1	french	french	39	
2	french	italian	30	
3	french	none	30	
4	italian	french	1	
5	italian	italian	19	
6	italian	none	11	
7	other	french	35	
8	other	italian	35	

**Contingency Table**

With Data  With Counts

First Variable: vine, music, n  
Second Variable: vine, music, n  
Counts: vine, music, n

Filter data  
 Insert results into data table

Buttons: Help, Cancel, OK

**Contingency Table with Counts -- Selected Fields: vin, tónlist, n**

	N	i	f	Total
c	43 (17.6955%)	35 (14.4033%)	35 (14.4033%)	113
a	30 (12.3457%)	30 (12.3457%)	39 (16.0494%)	99
b	11 (4.5267%)	19 (7.8189%)	1 (0.4115%)	31
Total	84	84	75	243

**Chi-Square**

df	4
value	18.2792
p Value	0.0011



# Ólíkar nálganir við túlkun

Núlltilgátan tilgreinir að engin tengsl séu í töflunni. Þetta má túlka á tvo vegu eftir því hvers eðlis rannsóknin er.

Engin tengsl birtast þannig að skilyrt dreifing verður eins fyrir öll gildi frumbreytunnar. Þessi túlkun hentar þegar líta má svo á að frumbreytan skilgreini ólík þýði.

Engin tengsl felur í sér að breyturnar eru innbyrðis óháðar. Þessi túlkun hentar ef líta má svo á að tvær breytur séu mældar í einu þýði.

## Tilraunaaðstæður

Í víndæminu er eðlilegt að líta svo á að uppruni tónlistar skilgreini þrjú ólík þýði, þ.e. vínbúð með franskri, ítalskri eða tónlist af öðrum uppruna. Gildi frumbreytunnar eru ákvörðuð. Því er eðlilegt að spyrja hvort hlutfallsleg skipting sé eins hjá hópnum.

## Mældar breytur

Kynferði og það hvort viðkomandi dettur í það eða ekki eru mældar breytur. Við getum litið á þetta sem þýði karla og þýði kvenna. En það er einnig eðlilegt að túlka þetta sem spurningu um það hvort breyturnar tvær séu háðar eða óháðar.

Önnur túlkunaraðferðin útilokar ekki endilega hina heldur eru þetta oft valkostir.

# Mátgæði (*goodness of fit*)

Með mátgæðum er kannað hversu vel dreifing passar við þá dreifingu sem við búumst við vegna forþekkingar eða fræðikenningar.

Við gefum upp raun- og fræðitíðni. Hér er fræðitíðnin miðuð við að slys dreifist jafnt á alla virka daga. Ég vel **Statistics/ Non-parametrics/Chi Squared**, tilgreini dálkana með raun- og væntitíðni og smelli á OK.

Niðurstaðan gefur til kynna frávik frá jöfnum líkindum sé ekki afgerandi.

The screenshot shows the CRUNCHIT! 2.0 interface. A data table is visible with columns for Row, Col 1, Col 2, Col 3, and Col 4. The data is as follows:

Row	Col 1	Col 2	Col 3	Col 4
#	Day	Count	Expected	
1	Wed	159	133.4	
2	Tue	126	133.4	
3	Fri	113	133.4	
4	Thu	136	133.4	
5	Mon	133	133.4	
6				

A dialog box titled "Chi Squared - Goodness of Fit" is open, showing "Observed Counts" and "Hypothesized Counts or Probabilities". The "Observed Counts" table has columns for Day, Count, and Expected. The "Hypothesized Counts or Probabilities" table has columns for Day, Count, and Expected. Below the tables are checkboxes for "Filter data" and "Insert results into data table".

A second dialog box titled "Chi Squared -- Selected Fields: Count, Expected" is open, showing a table of results:

Statistic	Result
Chi-squared Statistic	8.4948
df	4
p Value	0.0750

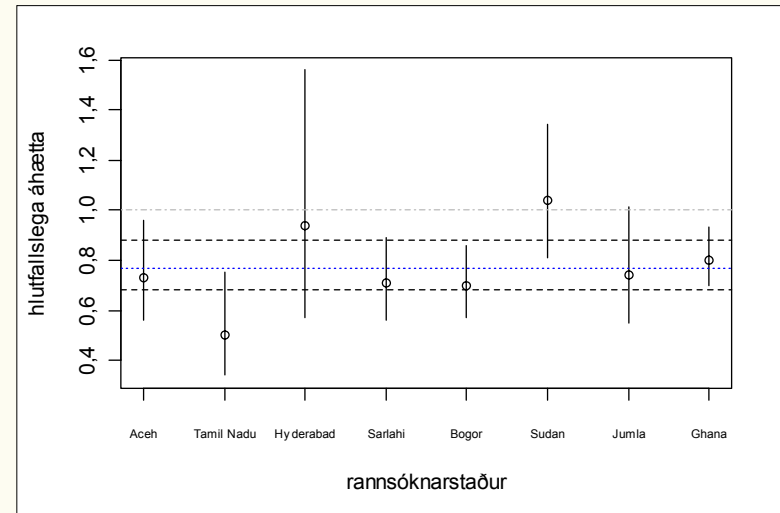
# Yfirgreining (*metaanalysis*)

Stök rannsókn gefur takmarkaðar upplýsingar. Ef hún hefur marktækt próf á tilgátu, gæti það verið vegna höfnunarmistaka; fastheldnimistök gætu einnig hafa orðið til þess að prófið varð ómarktækt.

Öryggisbil geta oft hjálpað með því að magnbinda óvissuna og tilgreina það talnabil sem sennilega innheldur þýðistöluna.

Nákvæmastar upplýsingar fást þó með yfirgreiningu, þ.e. samantekt á niðurstöðum margra rannsókna.

[http://www.unscn.org/layout/modules/resources/files/Policy\\_paper\\_No\\_13.pdf](http://www.unscn.org/layout/modules/resources/files/Policy_paper_No_13.pdf)



Hver stök rannsókn er villandi en allar átta teknar saman gefa öryggisbilið 0,68–0,88 miðað við 95% öryggi. Besta punktspá er 0,77, þ.e. 23% lækkun dánartíðni þegar vítamín A er gefið.

$$\text{Relative risk} = \frac{\text{hlutfall látinna í inngripi}}{\text{hlutfall látinna í samanburði}}$$