

# Visualizing data: the often neglected first step

Dr. Douwe Postmus (Unit Medical Statistics, Dept Epidemiology)

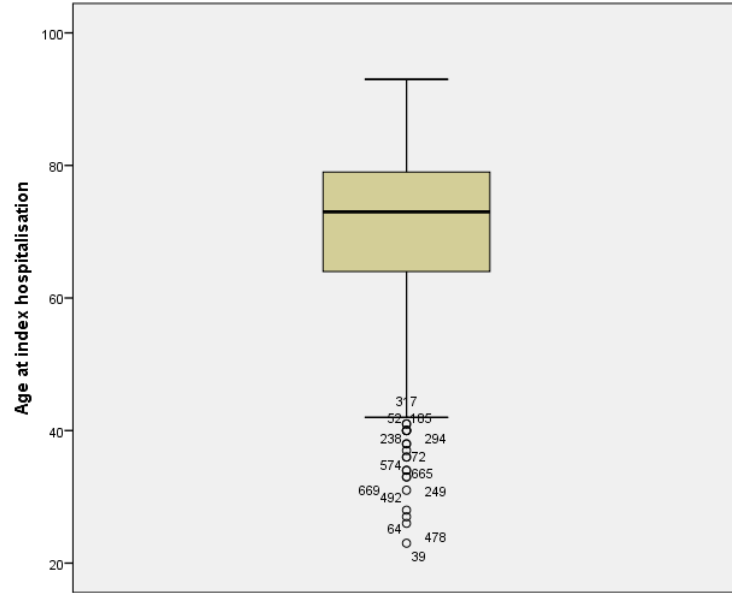
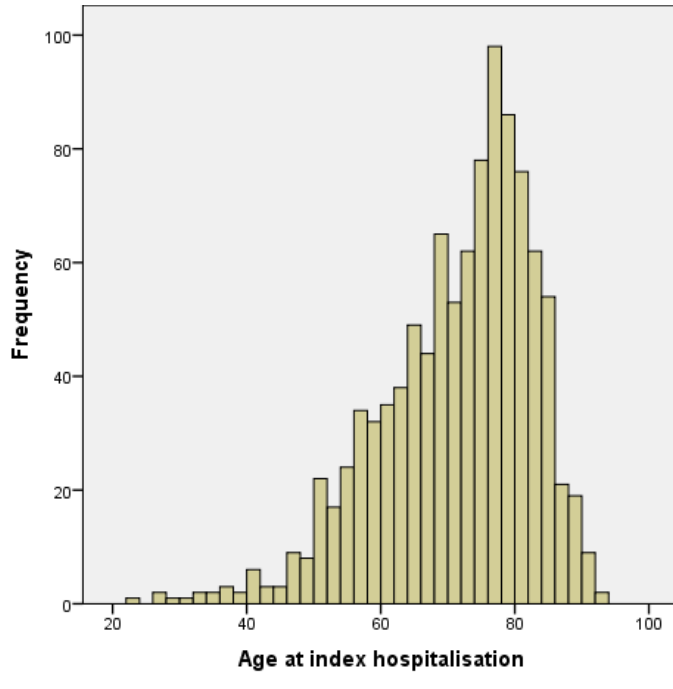


university of  
groningen



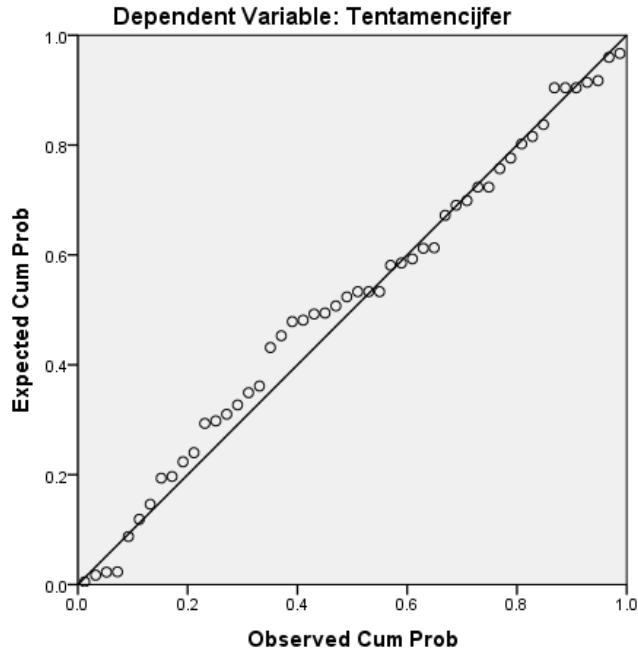
University Medical Center Groningen

# Data cleaning and validation

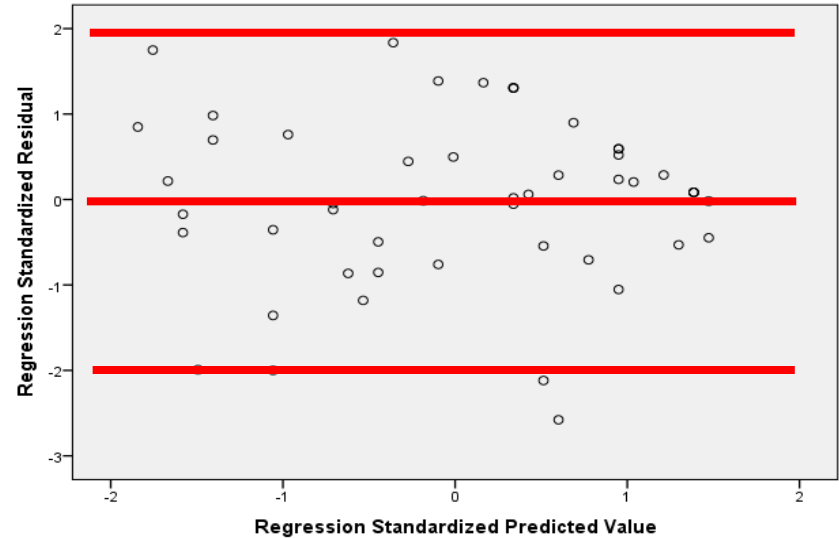


# Post-analysis: goodness of fit

Normal P-P Plot of Regression Standardized Residual

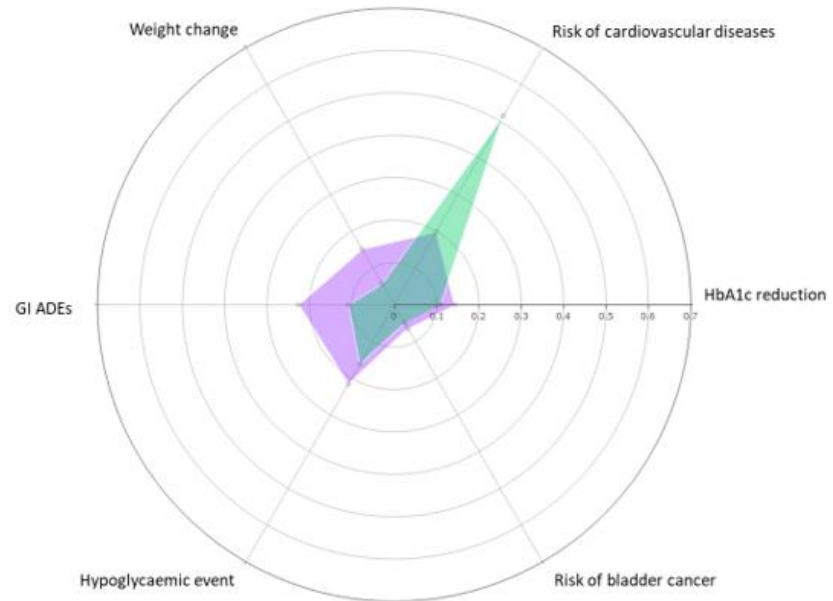


Scatterplot  
Dependent Variable: Tentamencijfer



# Post-analysis: results visualization

	Dapagliflozin n/N	Placebo Events/100 patient-years	Hazard Ratio (95% CI)	P Value for Interaction
<b>Primary outcome</b>				
eGFR decline $\geq 50\%$ , ESKD, or kidney or CV death				
Overall	197/2152	312/2152	4.6	7.5
Without CV disease	106/1339	175/1355	4.0	6.7
With CV disease	91/813	137/797	5.5	8.7
<b>Secondary outcomes</b>				
eGFR decline $\geq 50\%$ , ESKD, or kidney death				
Overall	142/2152	243/2152	3.3	5.8
Without CV disease	93/1339	154/1355	3.6	5.9
With CV disease	49/813	89/797	2.9	5.6
CV death or hospitalization for heart failure				
Overall	100/2152	138/2152	2.2	3.0
Without CV disease	24/1339	36/1355	0.8	1.3
With CV disease	76/813	102/797	4.3	6.1
All-cause death				
Overall	101/2152	146/2152	2.2	3.1
Without CV disease	33/1339	53/1355	1.1	1.8
With CV disease	68/813	93/797	3.8	5.4
<b>Prespecified exploratory CV outcomes</b>				
CV death, myocardial infarction, or stroke				
Overall	132/2152	143/2152	2.9	3.1
Without CV disease	41/1339	50/1355	1.4	1.7
With CV disease	91/813	93/797	5.2	5.5
First heart failure hospitalization				
Overall	37/2152	71/2152	0.8	1.6
Without CV disease	4/1339	13/1355	0.1	0.5
With CV disease	33/813	58/797	1.9	3.5
<b>Post-hoc exploratory CV/cardiorenal outcomes</b>				
CV death, myocardial infarction, stroke or heart failure hospitalization				
Overall	158/2152	195/2152	3.5	4.4
Without CV disease	44/1339	60/1355	1.5	2.1
With CV disease	114/813	135/797	6.6	8.3
All-cause death, myocardial infarction, stroke, heart failure hospitalization, or ESKD				
Overall	274/2152	376/2152	6.5	9.1
Without CV disease	118/1339	177/1355	4.5	6.8
With CV disease	156/813	199/797	9.6	13.1



university of  
 groningen



University Medical Center Groningen

# Today

- Use of graphical displays to summarize the patterns that are present in the data
- Helps with interpreting the results of a subsequent statistical test
- Supports statistical model building



# Example 1: How does a woman's behavior during pregnancy affect the infant's birth weight?

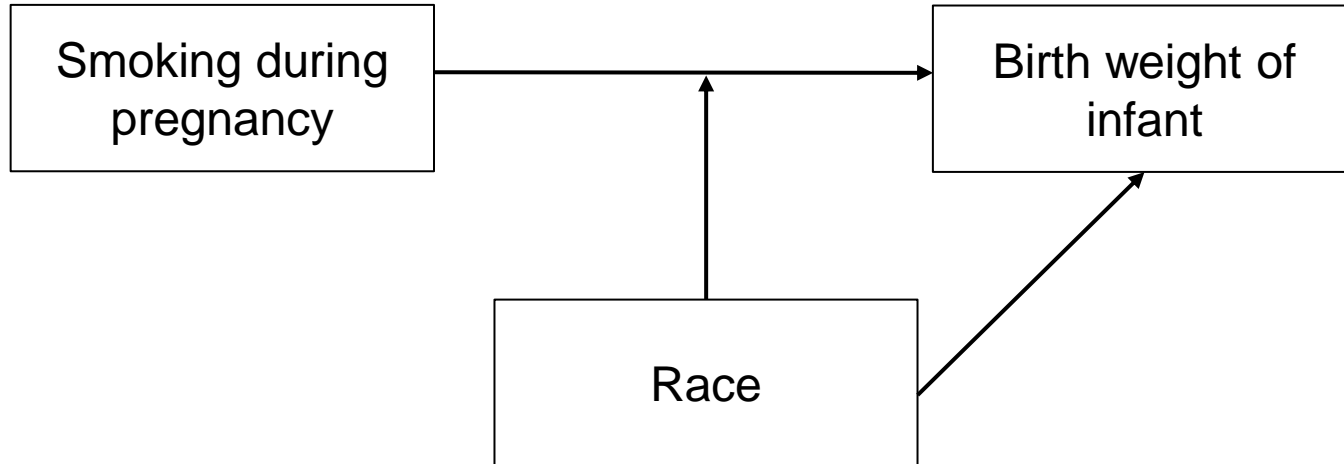
\* These data come from Appendix 1 of Hosmer and Lemeshow (1989), and were collected at Baystate Medical Center, Springfield MA, during 1986.

\* Low birth weight is an outcome that has been of concern to physicians for years. This is due to the fact that infant mortality rates and birth defect rates are very high for low birth weight babies. A woman's behavior during pregnancy (including diet, smoking habits, and receiving prenatal care) can greatly alter the chances of carrying the baby to term and, consequently, of delivering a baby of normal birth weight.

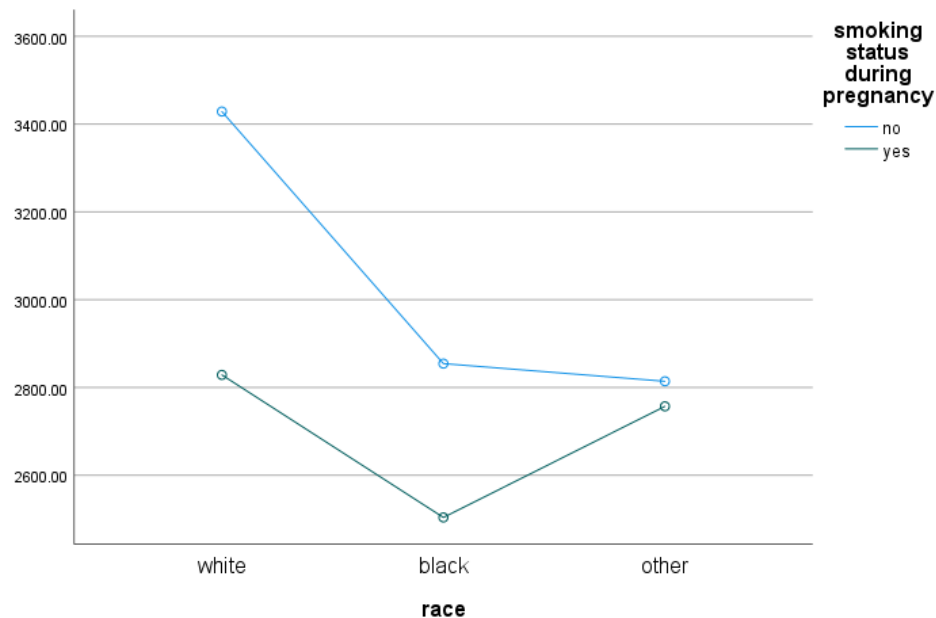
Columns	Variable	Abbreviation			
2-4	Identification Code	ID	40	Smoking Status During Pregnancy (1 = Yes, 0 = No)	SMOKE
10	Low Birth Weight (0 = Birth Weight $\geq$ 2500g, 1 = Birth Weight < 2500g)	LBW	48	History of Premature Labor (0 = None, 1 = One, etc.)	PTL
17-18	Age of the Mother in Years	AGE	55	History of Hypertension (1 = Yes, 0 = No)	HYPERT
23-25	Weight in Pounds at the Last Menstrual Period	LWT	61	Presence of Uterine Irritability (1 = Yes, 0 = No)	URIRR
32	Race (1 = White, 2 = Black, 3 = Other)	RACE	67	Number of Physician Visits During the First Trimester (0 = None, 1 = One, 2 = Two, etc.)	PVFT
			73-76	Birth Weight in Grams	BWT



# Conceptual model



# Interaction plot



In linear regression, the infant's mean birth weight is expressed as a linear function of the independent variables (regression equation)

Interaction plot: graphical display of the means for each combination of the levels of two categorical variables





# ANOVA table

## Tests of Between-Subjects Effects

Dependent Variable: birth weight in grams

Source	Type III Sum of Squares	df	Mean Square	F	Sig.
Corrected Model	14439242.57 <sup>a</sup>	5	2887848.514	6.183	<.001
Intercept	965464191.55	1	965464191.55	2066.969	<.001
race	5818289.482	2	2909144.741	6.228	.002
smoke	3318053.571	1	3318053.571	7.104	.008
race * smoke	2097537.495	2	1048768.747	2.245	.109
Error	85477810.075	183	467091.858		
Total	1738735950.0	189			
Corrected Total	99917052.646	188			

a. R Squared = .145 (Adjusted R Squared = .121)



# Post hoc tests for race

## Multiple Comparisons

Dependent Variable: birth weight in grams

Bonferroni

(I) race	(J) race	Mean Difference (I-J)	Std. Error	Sig.	95% Confidence Interval	
					Lower Bound	Upper Bound
white	black	384.0473*	151.09802	.036	18.9661	749.1285
	other	299.7247*	108.79826	.019	36.8476	562.6017
black	white	-384.0473*	151.09802	.036	-749.1285	-18.9661
	other	-84.3226	157.91324	1.000	-465.8707	297.2254
other	white	-299.7247*	108.79826	.019	-562.6017	-36.8476
	black	84.3226	157.91324	1.000	-297.2254	465.8707

Based on observed means.

The error term is Mean Square(Error) = 467091.858.

\*. The mean difference is significant at the 0.05 level.



# Example 2: predicting the 10-year risk of coronary heart disease (CHD)

556

Biometrical Journal 57 (2015) 4, 556–570 DOI: 10.1002/bimj.201300260

## Graphical assessment of incremental value of novel markers in prediction models: From statistical to decision analytical perspectives

Ewout W. Steyerberg<sup>\*1</sup>, Moniek M. Vedder<sup>1</sup>, Maarten J. G. Leening<sup>2,3</sup>, Douwe Postmus<sup>4</sup>, Ralph B. D'Agostino Sr.<sup>5</sup>, Ben Van Calster<sup>1,6</sup>, and Michael J. Pencina<sup>7</sup>

<sup>1</sup> Department of Public Health, Erasmus MC: University Medical Center Rotterdam, Rotterdam, The Netherlands

<sup>2</sup> Department of Epidemiology, Erasmus MC: University Medical Center Rotterdam, Rotterdam, The Netherlands

<sup>3</sup> Department of Cardiology, Erasmus MC: University Medical Center Rotterdam, Rotterdam, The Netherlands

<sup>4</sup> Department of Epidemiology, University Medical Center Groningen, University of Groningen, Groningen, The Netherlands

<sup>5</sup> Framingham Heart Study, Framingham, MA, USA

<sup>6</sup> Department of Development and Regeneration, KU Leuven, Leuven, Belgium

<sup>7</sup> Department of Biostatistics and Bioinformatics, Duke Clinical Research Institute, Duke University, Durham, NC, USA

Received 5 November 2013; revised 24 April 2014; accepted 25 May 2014

- Reference model: multivariable logistic model with sex, diabetes, and smoking as dichotomous predictors and age, systolic blood pressure, and total cholesterol as continuous predictors
- Does adding HDL cholesterol to an existing model improve risk prediction?
- Analysis based on 3264 participants from the Framingham Heart Study aged 30 – 74 years
- A total of 183 individuals developed CHD (5.6% 10-year cumulative incidence)

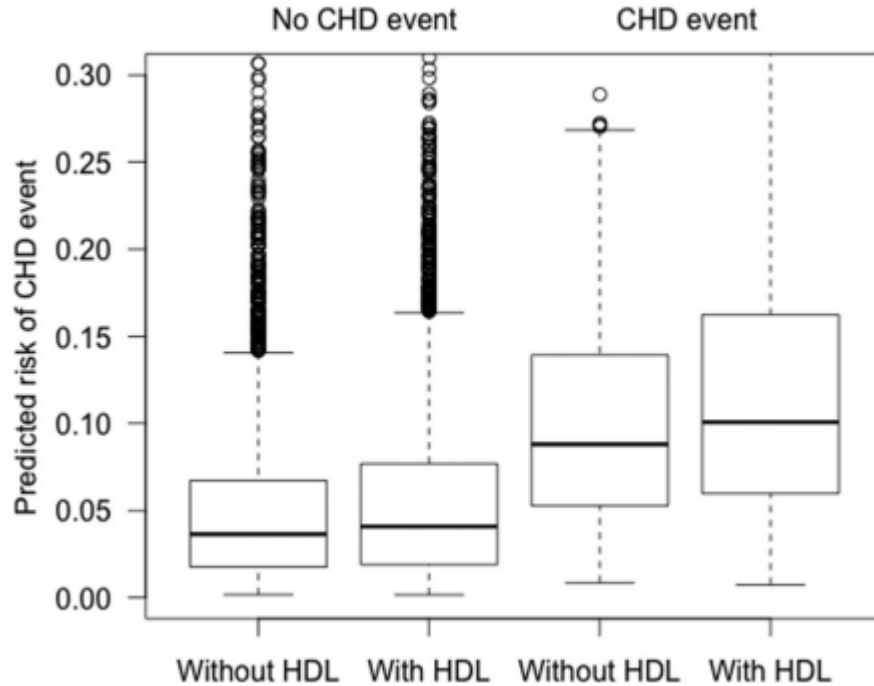


university of  
 groningen



University Medical Center Groningen

# Box plots stratified by CHD status - IDI



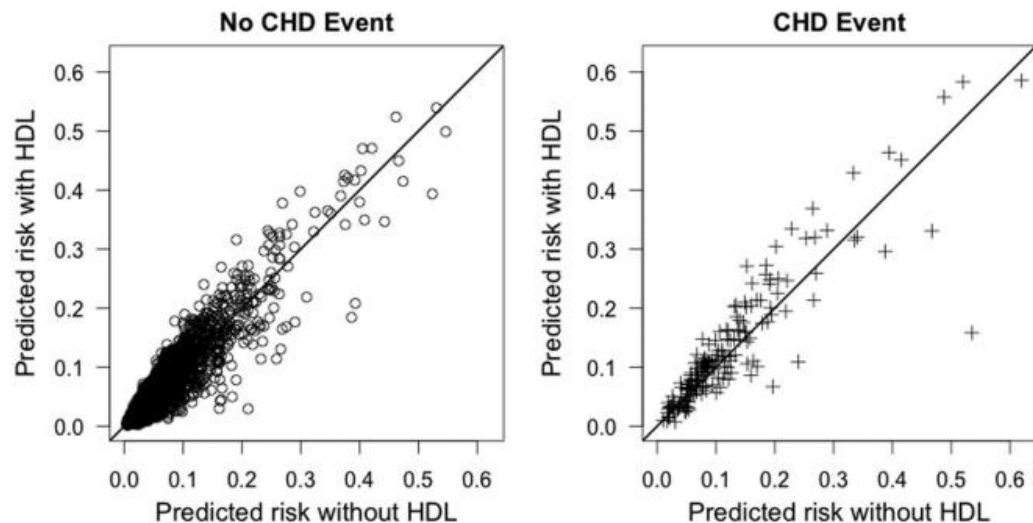
Discrimination slope =  
difference in mean  
predicted risks for those  
with and without the event

- Without HDL: 6.29%
- With HDL: 7.14%

Integrated discrimination  
index (IDI) = difference in  
discrimination slope  
 $7.14 - 6.29 = 0.85\%$



# Reclassification graphs - cNRI



Continuous net reclassification improvement (cNRI)

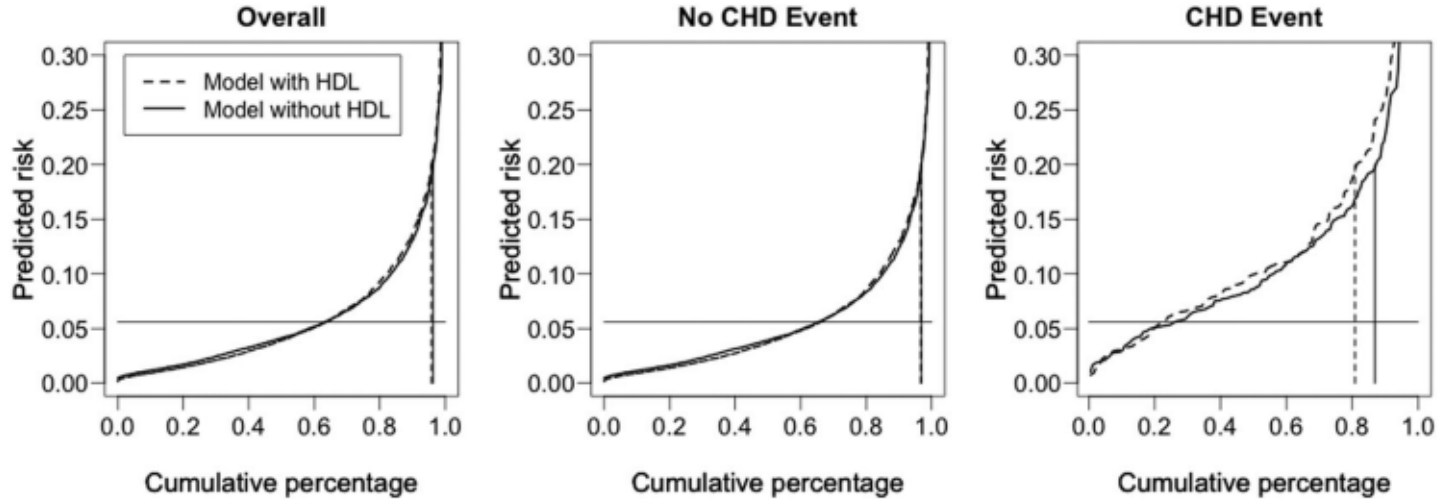
cNRI nonevents: a net 5.5% of nonevents receive lower predicted risks

cNRI events: a net 24.6% of those with events receive higher predicted risk

$$\text{cNRI} = 5.5\% + 24.6\% = 30.1\%$$



# Predictiveness curves – link between threshold and sensitivity/specificity

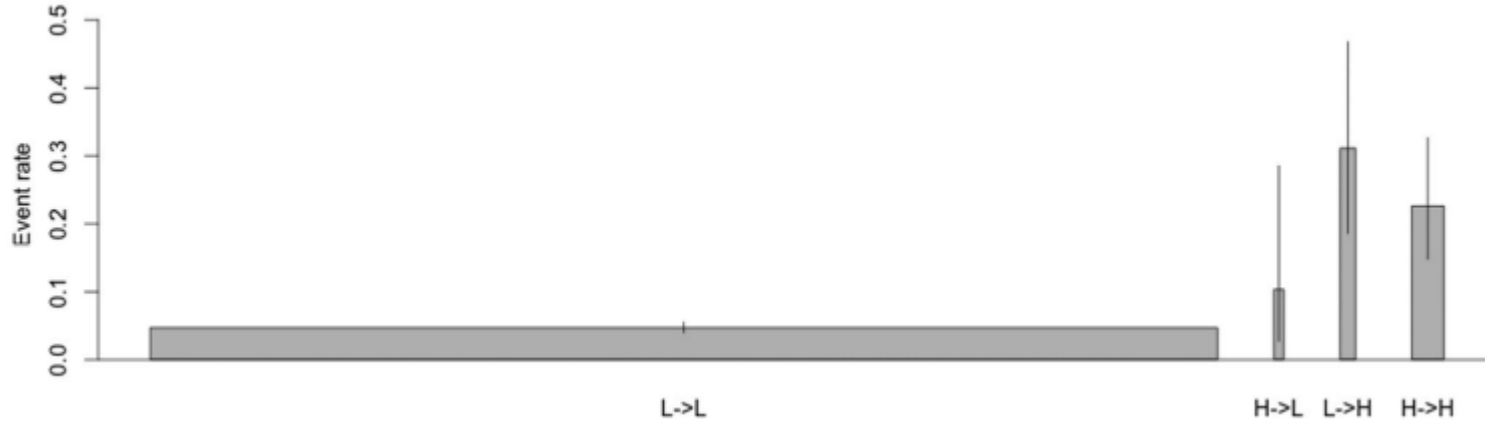


Specificity =  $P(- | \text{no CHD event}) = 96.82\%$  model with HDL vs  $96.66\%$  model without HDL

Sensitivity =  $P(+ | \text{no CHD event}) = 13.1\%$  model with HDL vs  $19.1\%$  model without HDL



# Net reclassification risk graph



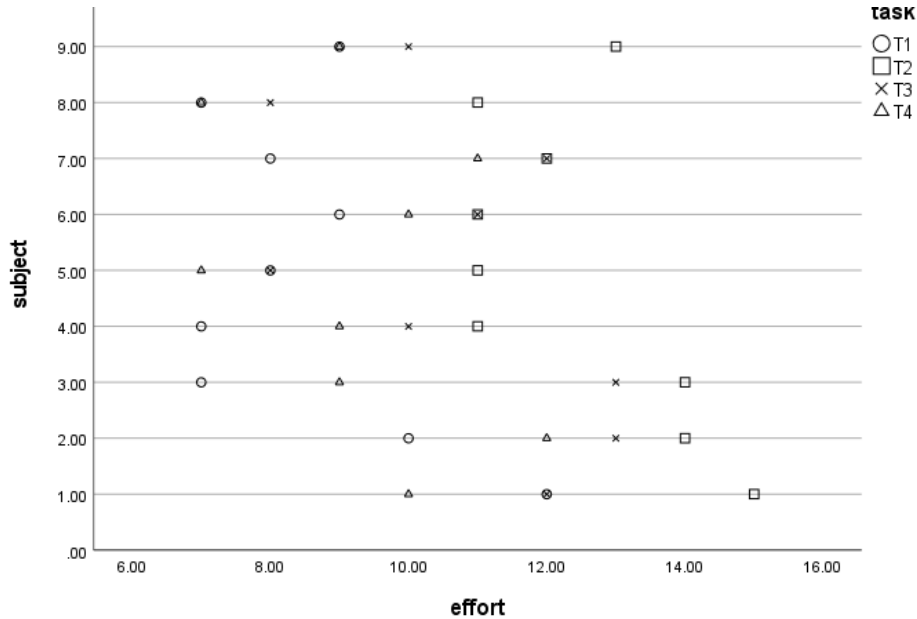
# Example 3: randomized block design

- Nine different subjects are asked to complete four different tasks
- The objectives are
  - To assess whether there are systematic differences in the complexity of the four tasks
  - To estimate the between-subject variability in task proficiency
- This experiment is an example of a randomized block design with task as a fixed effect and subject as a random effect





# Dot plot



Tasks 1 and 4 seem to take the least effort while task 2 seems to take the most effort

Moderate between-subject variability => interclass correlation (ICC)

# Results

Linear mixed-effects model fit by REML

Data: data

AIC	BIC	logLik
133.1308	141.9252	-60.56539

Random effects:

Formula: ~1 | Subject  
(Intercept) Residual  
StdDev: 1.332465 1.100295

Fixed effects: effort ~ Task

	Value	Std.Error	DF	t-value	p-value
(Intercept)	8.555556	0.5760123	24	14.853079	0.0000
TaskT2	3.888889	0.5186838	24	7.497610	0.0000
TaskT3	2.222222	0.5186838	24	4.284348	0.0003
TaskT4	0.666667	0.5186838	24	1.285304	0.2110

Correlation:

(Intr)	TaskT2	TaskT3	
TaskT2	-0.45		
TaskT3	-0.45	0.50	
TaskT4	-0.45	0.50	0.50

Standardized Within-Group Residuals:

Min	Q1	Med	Q3	Max
-1.80200345	-0.64316591	0.05783115	0.70099706	1.63142054

Number of Observations: 36

Number of Groups: 9

	numDF	denDF	F-value	p-value
(Intercept)	1	24	455.0075	<.0001
Task	3	24	22.3556	<.0001

$$ICC = 1.332465^2 / (1.332465^2 + 1.100295^2) = 0.59$$

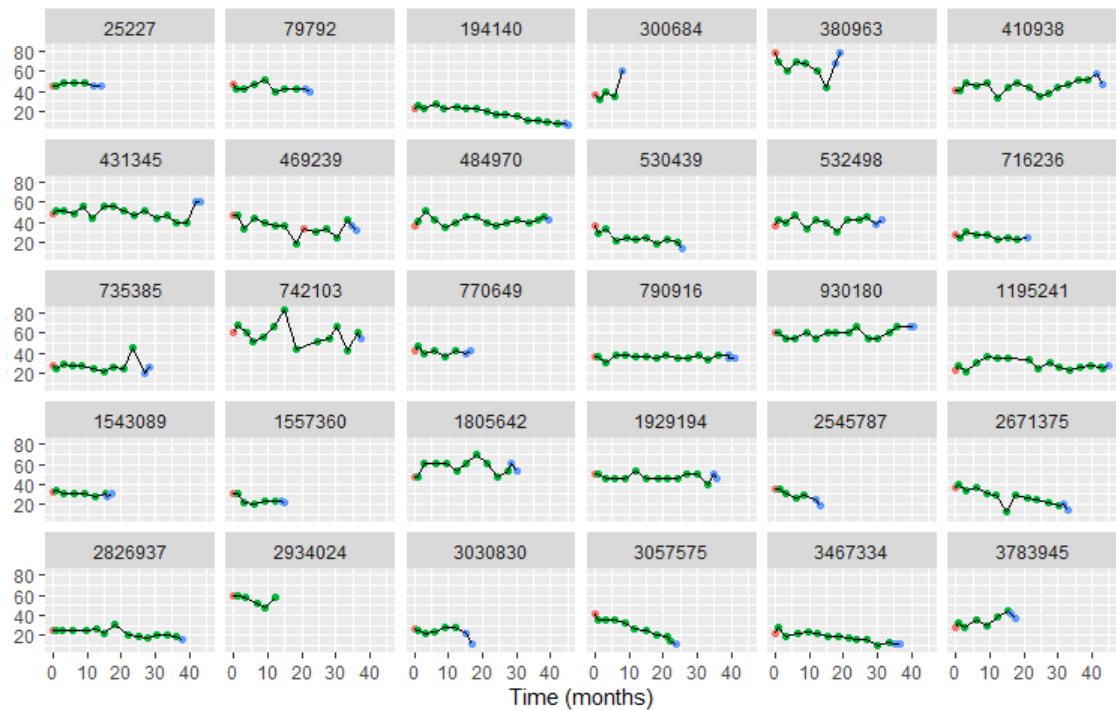


# Example 4: longitudinal biomarker measurements

- Objective: model the serial trends in a biomarker through a linear mixed effects model
- The fixed effect structure models the average biomarker trajectory
  - Linear (one-slope), piece-wise linear (two or more slopes), quadratic, cubic spline?
- The random effects model the subject-level deviations from the average trajectory
  - Which terms are needed to appropriately model the subject-level deviations?



# Individual profiles plot (trellis graph)



“An intelligent summary of data is often sufficient to fulfil the purposes for which the data were gathered, and more formal techniques such as confidence intervals and hypothesis tests sometimes add little to an investigator’s understanding”

John A. Rice, *mathematical statistics and data analysis, second edition*





[www.umcg.nl](http://www.umcg.nl)

contact: [d.postmus@umcg.nl](mailto:d.postmus@umcg.nl)



university of  
groningen



University Medical Center Groningen