**Slide 1**

# Large-scale meta-analysis of individual participant data: methods and challenges.

The Emerging Risk Factors Collaboration (ERFC)

Stephen Kaptoge
Department of Public Health and Primary Care
University of Cambridge
Cambridge, UK

Cambridge Statistics Discussion Group (CSDG) meeting, Cambridge (25 March 2013)

1

---

**Slide 2**

# Outline

- Meta-analysis
  - What, why, and how in general?

- Overview of IPD methods.
  - The Emerging Risk Factors Collaboration (ERFC).

- Examples of approaches to meta-analysis of:
  - Shape of dose-response relationships and specificity.
  - Adjusted main effects with correction for regression dilution.
  - Interactions (separating within vs between-study components).
  - Potential public health impact.
  - Risk prediction measures.

- Remarks
  - Stata programs made available for IPD meta-analysis.
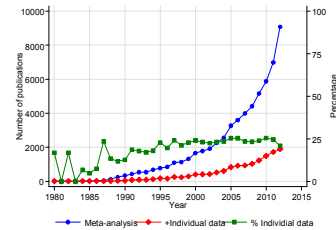
2

---

**Slide 3**

# Background

- An increasing number of molecular and genetic markers are being proposed as important predictors and/or causes of chronic diseases such as cardiovascular disease (CVD).

- Trend accelerated by the advent of technologies that enable rapid measurement of large numbers of blood proteins or genes.

- Reliable demonstration that a particular marker is relevant to CVD may have important implications for prediction and prevention (e.g. blood cholesterol).

- Meta-analysis *i.e. "use of statistical methods to combine results from individual studies"* should help improve statistical power to detect important associations and other features.
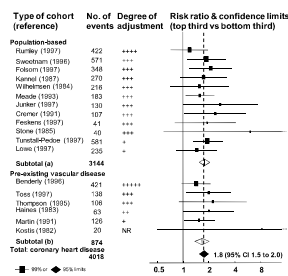
3

---

**Slide 4**

# Meta-analysis

- Use of meta-analysis has been growing in recent decades.
  - Literature-based meta-analysis.
  - Individual participant data (IPD) meta-analysis.



4

---

**Slide 5**

# Literature-based meta-analysis

**Prospective studies of fibrinogen and coronary heart disease by 1998**



Danesh et al, JAMA 1998

**Limitations:**
- Heterogeneity
  - Non-standardized definition of outcomes.
  - Inconsistent adjustment for confounding.
- Inability to assess shape of association.
- Impossible to assess effect modification by individual level characteristics.
- Associations not corrected for measurement error.

**Solution:**
- Large-scale meta-analysis using *individual participant data (IPD)*.

5

---

**Slide 6**

# Why IPD meta-analysis?

- Enables better standardised and more detailed analysis.
  - Standardised definition of outcomes and coding of variables.

- Improved power to assess dose-response relationships.

- Consistent adjustment for confounding across studies.

- Ability to correct estimates for regression dilution bias.

- Improved power to consistently assess subgroup effects.

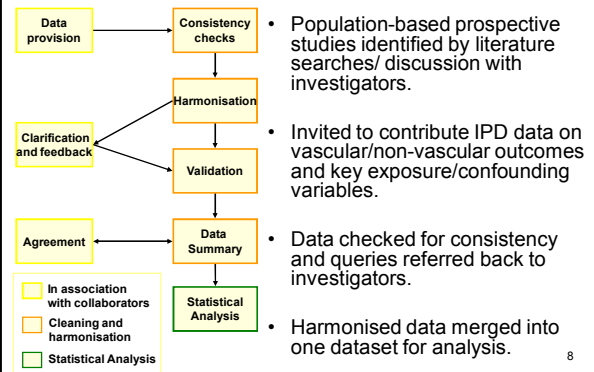- Evaluation of the utility of markers for risk prediction.

6

1

## The Emerging Risk Factors Collaboration (ERFC)

- A collation of IPD from 2.2 million participants in 143 prospective studies of cardiovascular disease (CVD) outcomes and cause-specific mortality.

- Aim: To quantify in detail the associations between established and emerging risk factors with CVD outcomes and cause-specific mortality, and to assess improvements in risk prediction.
  - Lipids and lipid-related markers (triglycerides, cholesterol, lipoprotein(a), apolipoproteins). JAMA 2009, JAMA 2012.
  - Inflammation markers (C-reactive protein, fibrinogen, leukocyte count, albumin). JAMA 2005, Lancet 2010, NEJM 2012.
  - Glycaemia markers (glucose, HbA1c) Lancet 2010.
  - Adiposity (BMI, waist, hip, waist:hip ratio). Lancet 2011.

7

## ERFC data collation process



- Population-based prospective studies identified by literature searches/ discussion with investigators.

- Invited to contribute IPD data on vascular/non-vascular outcomes and key exposure/confounding variables.

- Data checked for consistency and queries referred back to investigators.

- Harmonised data merged into one dataset for analysis.

8

## Study designs and methods in general

- Majority (>90%) are prospective cohort studies; a few clinical trials (mostly placebo arms); a few nested case-control studies for some novel markers; and one or two case-cohort studies.

- Principal analysis by 2-stage approach, first deriving study-specific estimates and then pooling by random effect meta-analysis.
  - Cox regression stratified by sex for cohort and clinical trials.
  - Conditional logistic regression for nested case-control studies.
  - Weighted Cox regression for case-cohort studies.

- Sensitivity analysis using fixed effect meta-analysis or 1-stage approach (i.e. stratified Cox model).e.g. for rare outcomes.

9

## Univariate meta-analysis (stage 1)

- In cohort studies, Cox-proportional hazards models fitted to each study.

$$\log(h_{ski}(t \mid E_{si}, X_{si})) = \log(h_{0sk}(t)) + \beta_s E_{si} + \gamma_s X_{si}$$

for each study $s = 1,...,S$, with strata $k = 1,...,K_s$, and individuals $i = 1,...,n_s$ with exposure of interest $E_{si}$ and other covariates $X_{si}$

- For most studies $K_s = 1$ or $2$, representing sex and in a few $K_s > 2$ in the presence of trial arm.

10

## Univariate meta-analysis (stage 2)

- $\beta_s$ = log hazard ratio per unit higher exposure in study $s$, adjusted for the confounding effects of the covariates $X_{si}$.

$$\hat{\beta}_s = \beta_s + \varepsilon_s; \text{ where } \varepsilon_s \sim N(0, v_s)$$
$$\beta_s = \beta + \eta_s; \text{ where } \eta_s \sim N(0, \tau^2)$$

- $\beta$ = average log hazard ratio, combining the within-study associations, while allowing for heterogeneity between studies (i.e. $\tau^2 > 0$).
  - DerSimonian & Laird moment estimator of $\tau^2$ is used, and the impact of heterogeneity is quantified using $I^2$ statistic.

- Stata program --metan-- for univariate meta-analysis.

11

## Multivariate meta-analysis

- At times interest is on meta-analysis of multiple correlated parameters, such as contrasts for categorical exposures, in which case multivariate meta-analysis is appropriate.

- Where now, $\mathbf{\beta}_s$ is a vector of log hazard ratios in study $s$ with known within study covariances $\Sigma_s$, adjusted for the covariates $X_{si}$.

$$\hat{\mathbf{\beta}}_s = \mathbf{\beta}_s + \mathbf{\varepsilon}_s; \text{ where } \mathbf{\varepsilon}_s \sim MVN(\mathbf{0}, \Sigma_s)$$
$$\mathbf{\beta}_s = \mathbf{\beta} + \mathbf{\eta}_s; \text{ where } \mathbf{\eta}_s \sim MVN(\mathbf{0}, \Omega)$$

- Stata program --mvmeta-- for multivariate meta-analysis.

--mvmeta-- White IR The Stata Journal 2009(9)40-56

12

## Regression dilution bias

- Noise in exposures and covariates (e.g. due to measurement errors or normal within person variability) often leads to regression dilution bias when using only a single baseline measurement in regression models.

$$X = Z + U; \text{ where } Z \text{ is true value, } U \sim N(0, \sigma_u^2)$$

$$Y = \beta_{0z} + \beta_z Z + \varepsilon; \text{ where } \varepsilon \sim N(0, \sigma_\varepsilon^2)$$
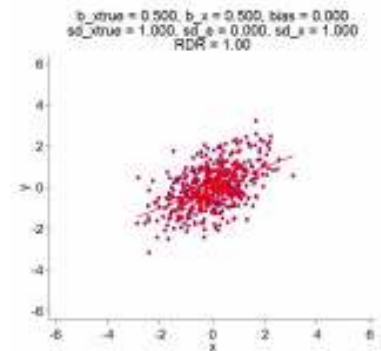
$$Y = \beta_{0x} + \beta_x X + \varepsilon; \text{ where } \varepsilon \sim N(0, \sigma_\varepsilon^2)$$

$$\beta_z = \frac{\beta_x}{\lambda} \text{ where } \lambda = \frac{\sigma_z^2}{\sigma_z^2 + \sigma_u^2}$$

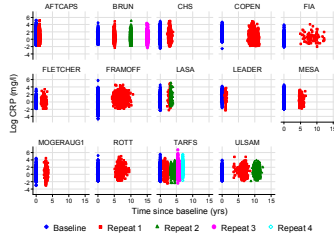- Rarely would the true value of X (i.e. Z) be known, but repeat measures of X can provide unbiased assessments.

## Regression dilution bias demo

## Correction for regression dilution

- Repeat measurements of exposure and other covariates allow for correction for impact of regression dilution bias simultaneously.



- Repeat measures are regressed on baseline measures to estimate regression dilution ratios (RDRs) over time.

## Regression dilution ratios (RDRs)

- Repeat measures are regressed on baseline measures to estimate regression dilution ratios (RDRs) over time.



RDR (95% CI)
R1: 0.63 (0.55, 0.71)
R2: 0.67 (0.61, 0.73)
R3: 0.48 (0.39, 0.57)

RDR (95% CI)
R1: 0.54 (0.47, 0.60)
R2: 0.54 (0.47, 0.60)
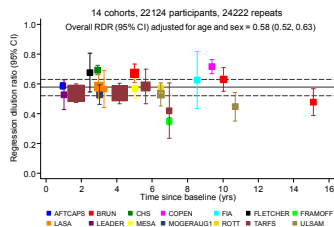R3: 0.58 (0.47, 0.70)
R4: 0.42 (0.23, 0.61)

- Linear mixed modelling used to fit a joint regression calibration model to data from all cohorts allowing for random effects across studies and individuals.

## RDRs and long-term "usual" levels

- RDRs of log CRP over time in 14 cohorts.



- Long-term average ("usual") levels of exposures predicted from the linear mixed regression calibration models and used as new exposures/covariates in the regression models for association with disease outcomes.
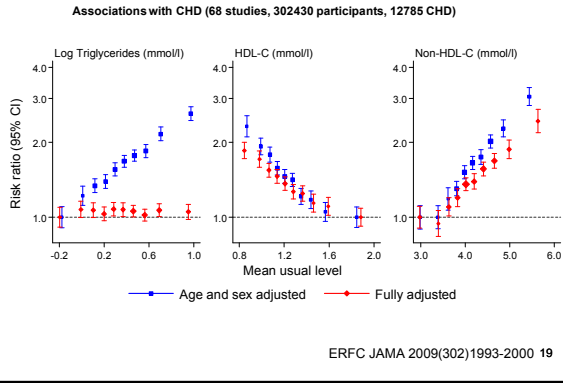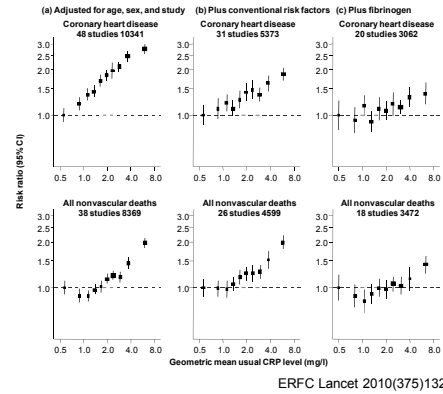
## Shape of dose-response relationships

- Categorise exposure into deciles (or appropriate bins) and calculate study-specific RR estimates vs. an appropriate reference group, adjusted for confounders as appropriate.

- Pool the study-specific log RRs by multivariate random effects meta-analysis (--mvmeta--, White IR) and plot against pooled mean exposure values within deciles.

- Sensitivity analysis by meta-analysis of continuous shape of association based on fractional polynomial models (--metacurve--, Royston P).
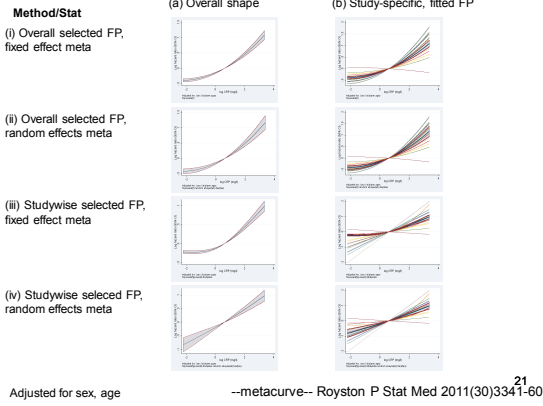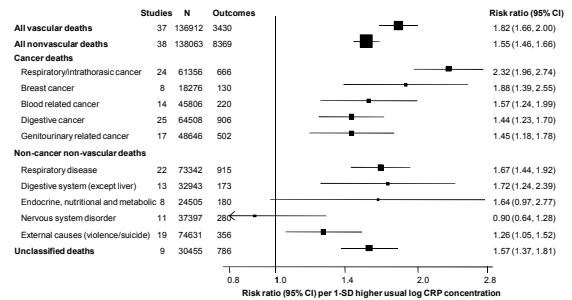
## Shape of dose-response relationships (Lipids)

Associations with CHD (68 studies, 302430 participants, 12785 CHD)
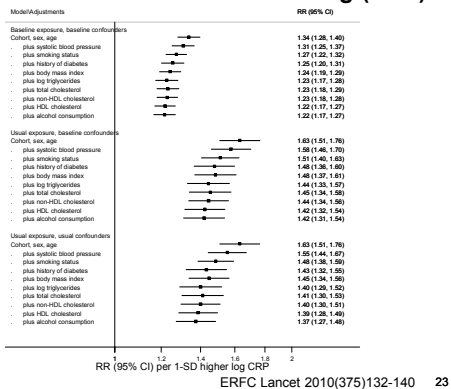


ERFC JAMA 2009(302)1993-2000 **19**

---

## Shape of dose-response relationships (CRP)



ERFC Lancet 2010(375)132-140 **20**

---

## Continuous shape of dose-response relationships (CRP)



Adjusted for sex, age

--metacurve-- Royston P Stat Med 2011(30)3341-60 **21**

---

## Specificity of associations (CRP)



ERFC Lancet 2010(375)132-140 **22**

---

## Regression dilution and confounding (CRP)



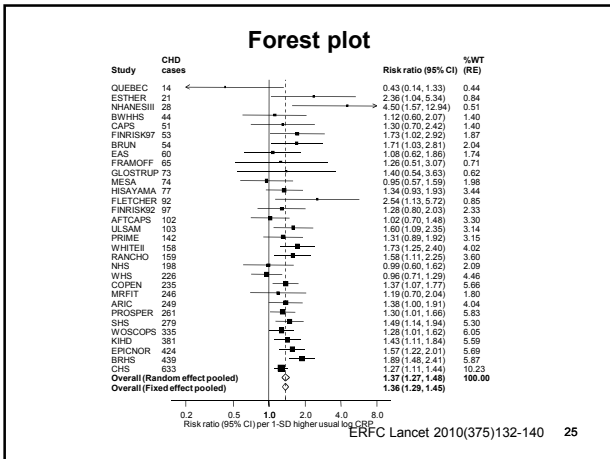ERFC Lancet 2010(375)132-140 **23**

---

## Confounding and between-study heterogeneity

RR for CHD per 1-SD higher usual log CRP adjusted for usual levels of confounders (31 studies, 91990 participants, 5373 CHD)

| | RR (95% CI) | Wald $\chi^2_1$ | $I^2$ (95% CI) |
|---|---|---|---|
| Adjusted for age, sex, and study | 1.63 (1.51, 1.76) | 149 | 51 (26, 68) |
| plus systolic blood pressure | 1.55 (1.44, 1.67) | 134 | 42 (12, 62) |
| plus smoking | 1.48 (1.38, 1.59) | 123 | 30 (0, 55) |
| plus history of diabetes | 1.43 (1.32, 1.55) | 78 | 42 (10, 62) |
| plus BMI | 1.45 (1.34, 1.56) | 85 | 31 (0, 56) |
| plus log$_e$ triglycerides | 1.40 (1.29, 1.52) | 64 | 35 (0, 58) |
| plus total cholesterol | 1.41 (1.30, 1.53) | 69 | 33 (0, 57) |
| plus non-HDL cholesterol[§] | 1.40 (1.30, 1.51) | 74 | 26 (0, 53) |
| plus cholesterol[§] | 1.39 (1.28, 1.49) | 71 | 25 (0, 52) |
| plus alcohol[§] | 1.37 (1.27, 1.48) | 65 | 26 (0, 53) |

ERFC Lancet 2010(375)132-140 **24**

## Forest plot



Risk ratio (95% CI) per 1-SD higher usual log CRP

ERFC Lancet 2010(375)132-140   **25**

---

## Joint effects (1)

- When effect modifiers are variables measured on individuals, e.g. age, BMI, etc interactions are most effectively assessed using within-study information.

$$\log(h_{ski}(t \mid E_{si}, X_{si})) = \log(h_{0sk}(t)) + \beta_s E_{si} + \gamma_s X_{si} + \delta_s E_{si} X_{si}$$

- Some potential effect modifiers are assessed only at the study level, e.g. assay methods, in which case assessment of interactions relies entirely on between-study comparisons (random-effects meta-regression).

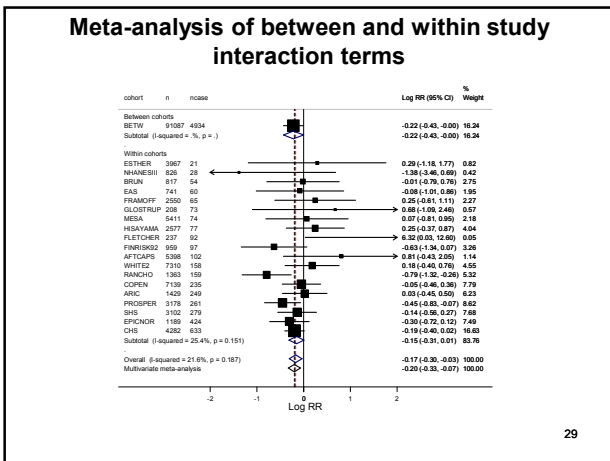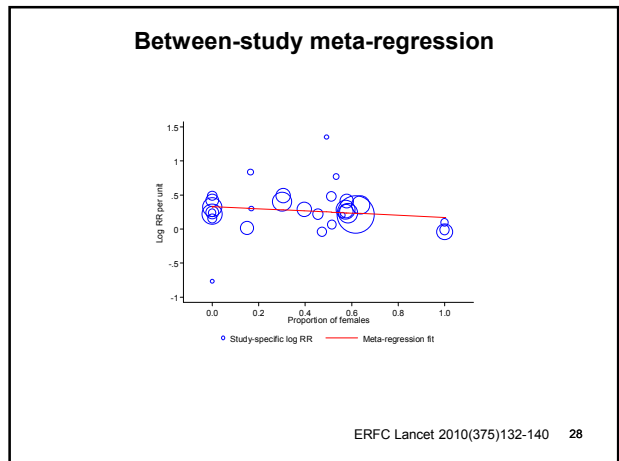$$\hat{\beta}_s = \beta_s + \varepsilon_s; \text{ where } \varepsilon_s \sim N(0, s_s^2)$$

$$\beta_s = \beta + \delta_B X_s + \eta_s; \text{ where } \eta_s \sim N(0, \tau^2)$$

**26**

---

## Joint effects (2)

- A few variables, notably sex and ethnic group, have potential interactions for which both within- and between- study information may be important.

- Studies of both sexes provide within-study information on sex interactions, while studies comprising one sex alone can only be used to assess interactions across studies.

- In ERFC we either:
  - Estimate both within- and between- study interactions and calculate inverse-variance weighted average.
  - Directly use multivariate meta-analysis to borrow information.

**27**

---

## Between-study meta-regression



ERFC Lancet 2010(375)132-140   **28**

---

## Meta-analysis of between and within study interaction terms



**29**

---

## Presentation of interactions

- For simple graphical representation of interactions of exposure and continuous effect modifiers, categorise the latter into tertiles and use multivariate meta-analysis to calculate the subgroup-specific estimates.

- But then …

- Base significance testing on the p-value of the pooled continuous interaction, which should be more powerful and less susceptible to artefacts from categorisation of the effect modifier.
  - Adopt more stringent p-values for significance of interactions.

**30**

## Joint effects with individual-level covariates

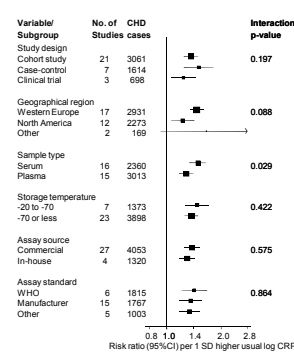| Variable/ Subgroup | Mean value | CHD cases | Interaction p-value |
|---|---|---|---|
| **Sex** | | | 0.015 |
| Male | | 3742 | |
| Female | | 1163 | |
| **Smoking history** | | | 0.710 |
| Other | | 3500 | |
| Current | | 1873 | |
| **Diabetes history** | | | 0.137 |
| Other | | 4663 | |
| Yes | | 710 | |
| **Age at survey (yrs)** | | | 0.022 |
| 40-59 | 52 | 2457 | |
| 60-69 | 64 | 1396 | |
| 70+ | 74 | 1215 | |
| **Systolic BP (mmHg)** | | | 0.453 |
| Bottom third | 117 | 1345 | |
| Middle third | 136 | 1789 | |
| Top third | 159 | 2239 | |
| **BMI (kg/m2)** | | | 0.969 |
| Bottom third | 23 | 1554 | |
| Middle third | 26 | 1813 | |
| Top third | 31 | 2006 | |

0.8  1.0  1.4  2.0  2.8
Risk ratio (95% CI) per 1-SD higher usual log CRP

ERFC Lancet 2010(375)132-140    **31**

---

## Joint effects with study-level covariates

| Variable/ Subgroup | No. of Studies | CHD cases | Interaction p-value |
|---|---|---|---|
| **Study design** | | | 0.197 |
| Cohort study | 21 | 3061 | |
| Case-control | 7 | 1614 | |
| Clinical trial | 3 | 698 | |
| **Geographical region** | | | 0.088 |
| Western Europe | 17 | 2931 | |
| North America | 12 | 2273 | |
| Other | 2 | 169 | |
| **Sample type** | | | 0.029 |
| Serum | 16 | 2360 | |
| Plasma | 15 | 3013 | |
| **Storage temperature** | | | 0.422 |
| -20 to -70 | 7 | 1373 | |
| -70 or less | 23 | 3898 | |
| **Assay source** | | | 0.575 |
| Commercial | 27 | 4053 | |
| In-house | 4 | 1320 | |
| **Assay standard** | | | 0.864 |
| WHO | 6 | 1815 | |
| Manufacturer | 15 | 1767 | |
| Other | 5 | 1003 | |

0.8  1.0  1.4  2.0  2.8
Risk ratio (95%CI) per 1 SD higher usual log CRP

ERFC Lancet 2010(375)132-140    **32**

---

## Proportional hazards (PH) assumption

- A key assumption in the Cox model is that of proportional hazards (PH), i.e. hazard ratios remain constant in time.

- We evaluate PH assumption in each study separately, by testing for interaction between the exposure and time.

$$\log(h_{ski}(t \mid E_{si}, X_{si})) = \log(h_{0sk}(t)) + \beta_s E_{si} + \boldsymbol{\gamma}_s X_{si} + \xi_s E_{si} t$$

- The non PH parameters $\xi_s$ are pooled across studies by random effect meta-analysis.

33

---

## PH assumption (alternative methods)

- Alternatively, can sum independent $\chi^2_1$ statistics from standard PH tests based on Schoenfeld residuals within each of $S$ studies, yielding a $\chi^2_S$ statistic testing H0 that PH holds in each study.

- Compare results with a 1-stage approach fitting a single model stratified by cohort and sex.

$$\log(h_{ski}(t \mid E_{si}, X_{si})) = \log(h_{0sk}(t)) + \beta_s E_{si} + \boldsymbol{\gamma}_s X_{si} + \xi E_{si} t$$

- where $\beta_s$ are separate fixed effects, and the focus is on the common estimate of $\xi$.
  - Can also assume that $\beta_s = \beta$, i.e. common effect across studies.

34

---

## Potential public health impact:
## Estimating years of Life Lost (YLL)

35

---

## Background

- Relative risks (RR) provide useful aetiological information on exposure (E) - disease (D) associations.

$$RR = \frac{p(D = 1 \mid E = 1, \mathbf{X} = \mathbf{x})}{p(D = 1 \mid E = 0, \mathbf{X} = \mathbf{x})}$$

- But alternative measures may be better suited for conveying the potential public health impact of a risk factor e.g.
  - Population attributable risk.
  - Short to medium term absolute risk (e.g. 10-year CVD risk).
  - Long term absolute risk (e.g. lifetime risk, or risk up to age 95y).
  - Years of life lost (YLL).

36

## Population attributable risk (PAR)

- By how much could population disease burden be reduced by eliminating a certain risk factor?

$$PAR = \left(1 - \frac{1}{RR}\right) \times p(E = 1 \mid D = 1)$$

- Lacks an individual perspective or interpretation.

- Ambiguous interpretation if exposure prevalence or if RRs importantly vary over time (e.g. by age).

37

## Short to medium term absolute risk

- What is the probability of disease (or death) occurring within t-years (e.g. 10-year CVD) given risk factor profile?

$$p(D = 1, T \leq t \mid E = 1, \mathbf{X} = \mathbf{x}) = 1 - S(t)$$
$$= 1 - S_0(t)^{\exp(\beta E + \boldsymbol{\beta}'\mathbf{x})}$$

- Has some individual perspective or interpretation.

- Interpretation is subject to length of time-horizon.

- Importance can be disproportionately weighted towards older age-groups, where incidence is higher.

38

## Long term absolute risk (e.g. lifetime risk)

- What is the probability of disease (or death) occurring over someone's lifetime conditional on survival to a certain age and risk factor profile (i.e. cumulative incidence, CI)?

$$p(D = 1 \mid age = age_{risk}, E = 1, \mathbf{X} = \mathbf{x}) = CI(age_{risk} \mid E = 1, \mathbf{X} = \mathbf{x})$$

- Individual perspective or interpretation.

- Should ideally take into account competing risks (e.g. death from non-cardiovascular disease).

- Gives more weight to events occurring at younger ages.

39

## Years of life lost (YLL)

- Conditional on survival to a certain age, how many years of life are lost because someone's risk factor profile promotes premature death?

- Essentially a difference in area under survival curves for exposed vs. unexposed from $age_{risk}$ to $age_{max}$.

$$YLL = \int_{age_{risk}}^{age_{max}} S(t \mid E = 1, \mathbf{X} = \mathbf{x}) dt - \int_{age_{risk}}^{age_{max}} S(t \mid E = 0, \mathbf{X} = \mathbf{x}) dt$$

- Gives a more intuitive perspective for the individual since it is expressed in life-year units rather than probability or RR scale.
  - Also it is calculated over the potential life-course of an individual.

40

## YLL estimation: Statistical inputs needed

- Age-at-risk specific hazard ratios (HRs) for exposure association with all-cause and cause-specific death (e.g. from ERFC).

- Population cause-specific and overall death rates by age and sex (e.g. EU 2000 death rates, in 5-year age groups).
  - Preferably for unexposed group (e.g. non-diabetics or non-smokers), but a work-around is possible if only available for total population.

- A computer with large memory, since estimating age-at-risk specific HRs involves expanding the dataset to obtain a record for each 5-year age-at-risk group.

41

## YLL estimation: Method in general

- Estimate sex and age-at-risk specific HRs for association of exposure with causes of interest (e.g. using ERFC data).

- Infer the expected death rates among the exposed (e.g. diabetics) by multiplying the reference population age-specific death rates by the estimated age-specific HRs.

- Derive the expected population survival curves for the exposed and unexposed groups, assuming exponential survival within each 5-year age-at-risk group.

- Estimate remaining life-years for each group by integrating both survival curves from each $age_{risk}$ to $age_{max}$ (cubic splines approx) and then calculate YLL as the difference.

42

7

## Slide 43

ERFC example 1: YLL lost due to diabetes

**43**

## Slide 44

### Data setup for modelling age-at-risk specific HRs …

```
+----------------------------------------------------------------------------------------------------------+
| idno    sex  ages   ageout  hxdiab  duration  ep_dead  agerisk  outage  agegrp5  _d   _t0    _t   _trisk |
+----------------------------------------------------------------------------------------------------------+
| C000074  1  59.99  70.3    0       10.3      0        59.99    60      55       0    0     .011   .011 |
| C000074  1  59.99  70.3    0       10.3      0        60       65      60       0    .011  5.01   5    |
| C000074  1  59.99  70.3    0       10.3      0        65       70      65       0    5.01  10     5    |
| C000074  1  59.99  70.3    0       10.3      0        70       70.3    70       0    10    10.3   .297 |
+----------------------------------------------------------------------------------------------------------+
| C000113  1  64.2   75.86   0       11.7      1        64.2     65      60       0    0     .8     .8   |
| C000113  1  64.2   75.86   0       11.7      1        65       70      65       0    .8    5.8    5    |
| C000113  1  64.2   75.86   0       11.7      1        70       75      70       0    5.8   10.8   5    |
| C000113  1  64.2   75.86   0       11.7      1        75       75.86   75       1    10.8  11.7   .863 |
+----------------------------------------------------------------------------------------------------------+
| C000135  1  49.35  61.93   1       12.6      0        49.35    50      45       0    0     .648   .648 |
| C000135  1  49.35  61.93   1       12.6      0        50       55      50       0    .648  5.65   5    |
| C000135  1  49.35  61.93   1       12.6      0        55       60      55       0    5.65  10.6   5    |
| C000135  1  49.35  61.93   1       12.6      0        60       61.93   60       0    10.6  12.6   1.93 |
+----------------------------------------------------------------------------------------------------------+
| C000217  1  58.65  71.16   1       12.5      1        58.65    60      55       0    0     1.35   1.35 |
| C000217  1  58.65  71.16   1       12.5      1        60       65      60       0    1.35  6.35    5    |
| C000217  1  58.65  71.16   1       12.5      1        65       70      65       0    6.35  11.4    5    |
| C000217  1  58.65  71.16   1       12.5      1        70       71.16   70       1    11.4  12.5   1.16 |
+----------------------------------------------------------------------------------------------------------+
```

Sex-specific Cox model: $s$ = study, $i$ = individual, $t$ = time in study

$$\log(h_{si}(t)) = \log(h_{si0}(t)) + \beta_1 diab_{si} + \beta_2 agerisk_{si} + \beta_3 agerisk_{si}^2 + \beta_4 diab_{si} \times agerisk_{si} + \beta_5 diab_{si} \times agerisk_{si}^2$$

**44**

## Slide 45

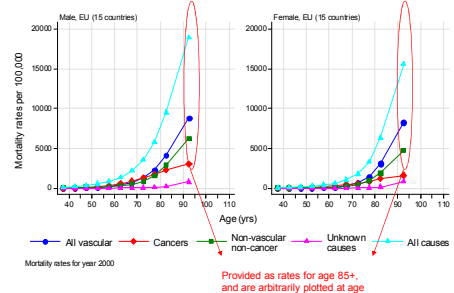### Input 1: Age specific HRs for cause-specific death (ERFC)



Inference from Cox model with diabetes status interacted with linear and quadratic terms for age-at-risk.

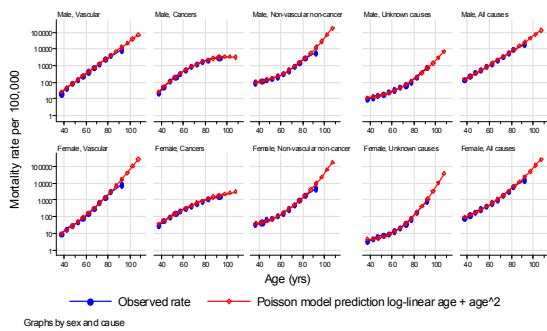**45**

## Slide 46

### Input 2: Population death rates (EU 2000, 15 countries)



Mortality rates for year 2000

Provided as rates for age 85+, and are arbitrarily plotted at age 92.5 y
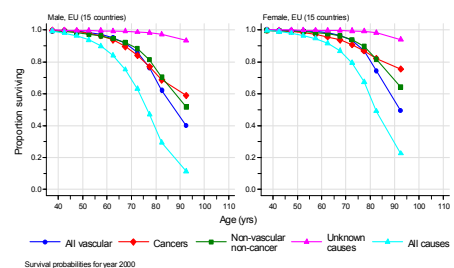
**46**

## Slide 47

### Input 2: Poisson smoothing/extrapolation of death rates



Graphs by sex and cause

**47**

## Slide 48

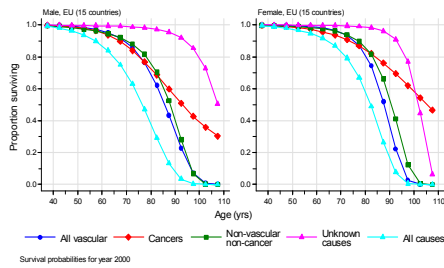### Input 3: Estimated population survival curves (1)



Survival probabilities for year 2000

$$p(survival \mid age_{risk} \geq 35) = \prod_{age_{risk}}^{age_{max}} e^{-5 \times \frac{IR_i}{100000}}$$

**48**

**Input 3: Estimated population survival curves (2)**

Male, EU (15 countries)  Female, EU (15 countries)

Proportion surviving — Age (yrs)

All vascular — Cancers — Non-vascular non-cancer — Unknown causes — All causes

Survival probabilities for year 2000

$$p(survival \mid age_{risk} \geq 35) = \prod_{age_{risk}}^{age_{max}} e^{-5 \times \frac{IR_i}{100000}}$$

49

---

**YLL Task 1: Survival by diabetes status**

Male, All causes  Female, All causes

Survival probability — Age (yrs)

No Diabetes — Diabetes

Graphs by sex and cause

50

---

**YLL Task 2: Remaining life years by diabetes status …**

Male, All causes  Female, All causes

Years of life remaining (yrs) — Age(yrs)

No Diabetes — Diabetes

Graphs by sex and cause

51

---

**YLL Task 3: Years of life lost with diabetes …**

Male, All causes  Female, All causes

Years of life lost (yrs) — Age (yrs)

Graphs by sex and cause

52

---

## Proportioning YLL by cause of death

- Calculation of YLL by cause of death requires additional constraints on the rates of 'other causes' of death when applying the cause-specific HRs to the death rates.
  - Exposure effect on the cause of interest is calculated while maintaining the death rates for all 'other causes' of death at their original death rates.

- Ensures that the sum of cause-specific YLLs is close to the YLLs directly estimated based on the all-cause mortality endpoint.

- For each age and sex group, denote the death rates as r_A (all-cause), r_V (vascular), r_C (cancer), r_N (non-vascular, non-cancer), and r_O (other), and the hazard ratios as HR_* with the same suffices.
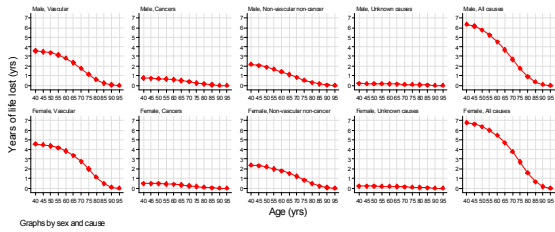
53

---

## Proportioning YLL by cause of death

- The rates add up so that r_A = r_V + r_C + r_N + r_O.

- For YLLs due to vascular death, compare survival curves derived from (r_V + r_C + r_N + r_O) to that derived from (HR_V*r_V + r_C + r_N + r_O).

- For YLLs due to cancer death, compare survival curves derived from (r_V + r_C + r_N + r_O) to that derived from (r_V + HR_C*r_C + r_N + r_O).

- For YLLs due to non-vascular, non-cancer death, compare survival curves derived from (r_V + r_C + r_N + r_O) to that derived from (r_V + r_C + HR_C*r_N + r_O), etc.
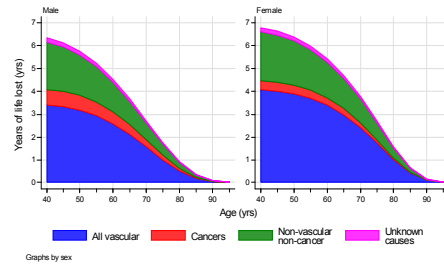
54

9

## YLL Task 4: YLL by cause and age-at-risk …



Graphs by sex and cause

**55**

## YLL Task 5: Proportionate YLL with diabetes by cause of death



Graphs by sex

**56**

---

ERFC example 2: YLL lost due to smoking

**57**

---

## Inferring death rates for the reference group

- The population death rates would rarely be available separately for the exposure groups of interest e.g. never smokers (0) – reference group, former smokers (1), and current smokers (2).

- To get death rates for never smokers (reference group), assume the overall population age-specific incidence rate is a weighted average

  $$IR = p_0*IR_0 + p_1*IR_1 + p_2*IR_2$$

  Given relative risks versus never smokers are

  $$RR_{10} = IR_1/IR_0 \text{ and } RR_{20} = IR_2/IR_0$$

- Can rewrite above equation as

  $$IR = p_0*IR_0 + p_1*RR_{10}*IR_0 + p_2*RR_{20}*IR_0$$

  from which

  $$IR_0 = IR/(p_0 + p_1*RR_{10} + p_2*RR_{20})$$

- Hence need age-specific prevalence of never ($p_0$), former ($p_1$), and current ($p_2$) smoking. Moreover, as smoking prevalence has declined over the years, period specific prevalence estimates are appropriate.
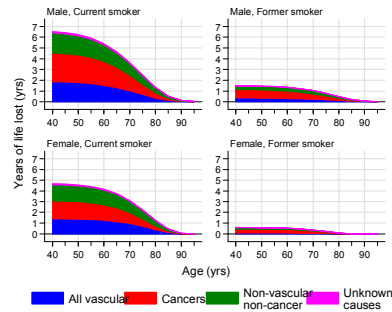
**58**

---

## Sex- and age-specific prevalence of smoking by decade



Inference from sex-specific multinomial logit model with linear and quadratic terms for age interacted with decade of survey.

**59**

## Proportionate YLL due to smoking by cause of death



**60**

## Remarks

- The outlined methods for calculation of YLL provide a useful way of combining overall population mortality rates and exposure association statistics from population-based studies to express the public health impact of risk factors.

- Future work would involve fully taking into account the impact of competing risks (e.g. cumulative incidence calculation) and finer subdivision by cause of death (e.g. lung v. other cancers for smoking).

- The results provide credible evidence that diabetes and current smoking are each associated with significantly shortened life expectancy. The impact of diabetes is greater in females than males, and the impact of smoking is greater in males.

61

## Risk prediction methods in ERFC

### (If time allows)

62

## Risk prediction - Background

- In addition to summarising the association between a risk marker and disease, there is interest in knowing the predictive ability that can be attributed to a marker when it is included in a model to predict future risk of CVD.

- Need to assess the predictive ability of risk models with and without inclusion of the marker of interest.

- In ERFC we have assessed:
  – Discrimination
  – Calibration
  – Reclassification

63

## The risk prediction model

$$\hat{S}_{ki}(t \mid \mathbf{x_i}, k) = \hat{S}_{0k}(t)^{\exp(\hat{\boldsymbol{\beta}}\mathbf{x}_i)}$$

- Gives the probability of surviving event free to at least time $t$ given risk factors $\boldsymbol{x}_i$.

- Interested in assessing the accuracy of the predictions made for all participants, given risk factors.

- Also interested in assessing the improvements in risk predictions when a new risk marker is added to some existing set of predictors.

64

## Discrimination

- The ability of a model to discriminate between different levels of risk.

- Assessed using:
  - Concordance statistics (C-index, AUROC)
  - The D measure or $R^2_D$

- The improvement in risk prediction upon addition of a new marker is quantified by assessing the change in discrimination measures e.g.
  C-index for new model – C-index for old model

65

## C-index for survival data

- Estimates the probability that the predicted order of failure is correct for a randomly selected pair of individuals.

- Calculation of the C-index (within a single study):

Count pairs of participants for which $t_i$ and $\hat{\beta}x_i$ are concordant ($n_c$), discordant ($n_d$) and undecided ($n_u$)
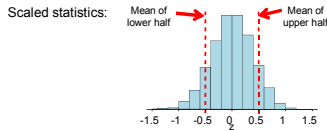
$$C = \frac{n_c + 0.5 n_u}{n_c + n_d + n_u}$$

Where $t_i$ = time in the study for participant $i$

$\hat{\beta}x_i$ = prognostic index/risk prediction for participant $i$

66

11

## D-statistic

- Measures the spread of observed risk across model predictions. Expressed as log-hazard ratio equivalent to comparing risk of outcome in the upper vs. lower half of the predicted risk distribution.

- Estimation steps
    1. Fit prediction model (eg: Cox PH)
    2. Transform linear predictor ($\beta X$) to scaled N(0,1) statistics (z)
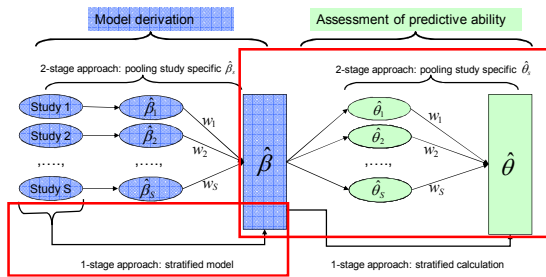    3. Fit outcome model (eg: Cox PH) to the scaled statistics, from which D is the coefficient of z

Scaled statistics: Mean of lower half    Mean of upper half

-1.5  -1  -0.5  0  0.5  1  1.5

67

---

## Risk prediction in ERFC

- Aim: To assess improvements in risk prediction by meta-analysis of data from observational studies.

- Measures used:
    - Discrimination: C-index, D measure, $R^2_D$
    - Reclassification: NRI, IDI

68

---

## ERFC prediction methods schemata



$W_i$ = study-specific weights for meta-analysis

Primary ERFC prediction analyses

69

---

## ERFC CRP-CHD prediction examples
### (37 studies, 165856 participants, 8806 CHD)

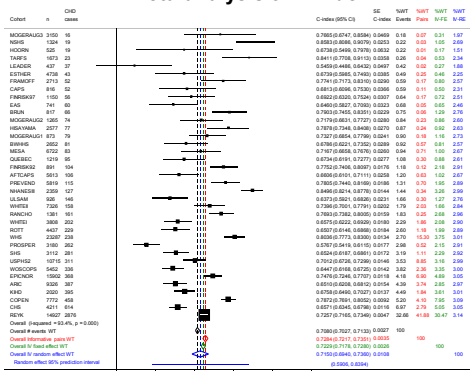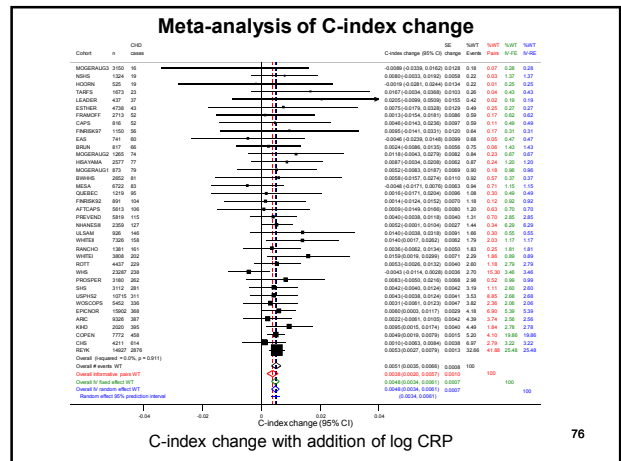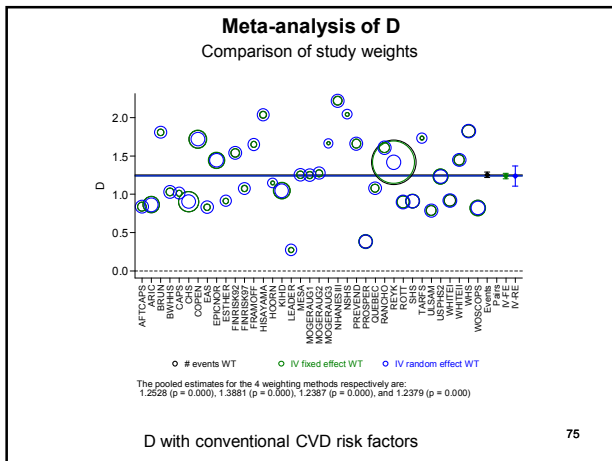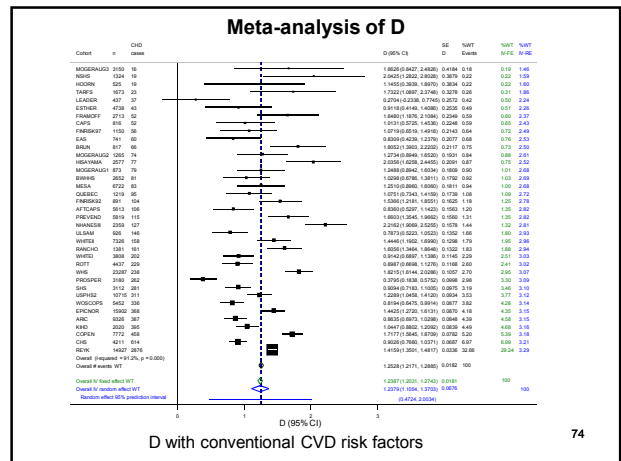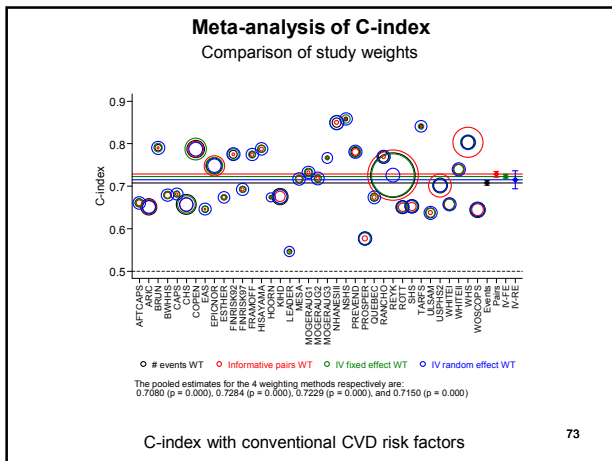| | | Multivariable adjusted log HR (SE) per 1-SD higher value or in comparison with reference category* | | | Heterogeneity |
|---|---|---|---|---|---|
| | Mean (SD) or n (%) | 1-stage Stratified† | 2-stage FE$ | 2-stage RE$ | I² (95% CI) |
| Age at survey (yrs) | 64.2 (8.6) | 0.567 (0.013) | 0.565 (0.013) | 0.529 (0.043) | 76 (67, 82) |
| Male sex | 81732 (49%) | NA | NA | NA | NA |
| Current smoking* | 35577 (21%) | 0.516 (0.024) | 0.529 (0.024) | 0.515 (0.050) | 63 (48, 73) |
| SBP (mmHg) | 131 (19) | 0.202 (0.009) | 0.203 (0.009) | 0.211 (0.017) | 30 (0, 53) |
| History of diabetes* | 10790 (7%) | 0.557 (0.038) | 0.587 (0.037) | 0.600 (0.049) | 24 (0, 49) |
| TCHOL (mmol/l) | 5.84 (1.06) | 0.234 (0.010) | 0.235 (0.010) | 0.216 (0.018) | 32 (0, 54) |
| HDL-C (mmol/l) | 1.27 (0.38) | -0.247 (0.014) | -0.240 (0.014) | -0.232 (0.023) | 52 (32, 67) |
| Log CRP (mg/l) | 0.55 (1.09) | 0.206 (0.012) | 0.207 (0.012) | 0.201 (0.013) | 9 (0, 38) |

70

---

## ERFC methods: Discrimination

- Discrimination (2-stage): C-index, D, $R^2_D$.
    - Stratified (1-stage) Cox model used to derive a single risk prediction model across studies.

    - Discrimination measures calculated within each study, including their changes therein.

    - Absolute values and changes combined by meta-analysis, <u>weighting by number of events</u>.

    - Sensitivity analyses using 2-stage derivation of the prediction model and also <u>alternative weights</u> for the meta-analysis of discrimination measures.
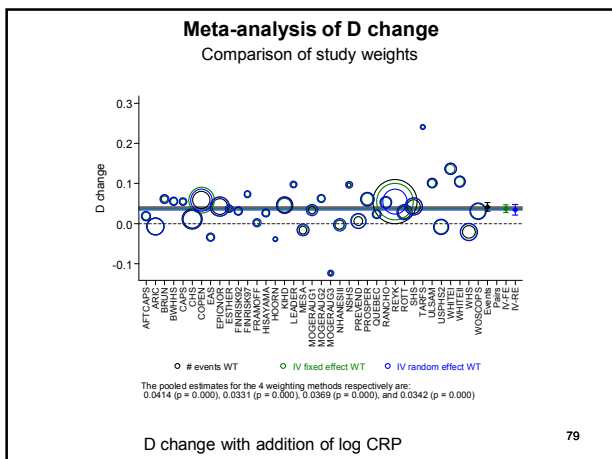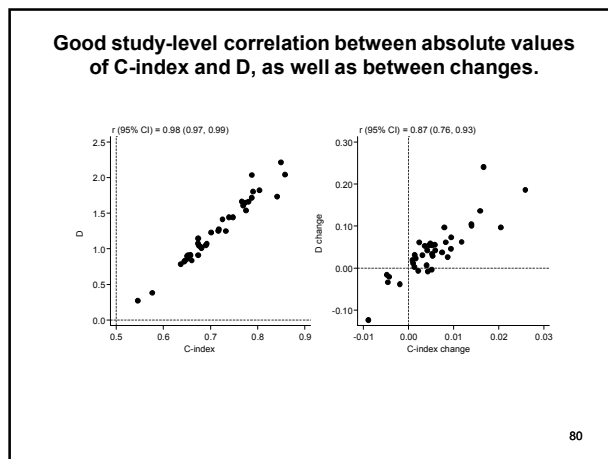
71

---

## Meta-analysis of C-index



C-index with conventional CVD risk factors

72

---

12

**Meta-analysis of C-index**
Comparison of study weights

C-index with conventional CVD risk factors

73



**Meta-analysis of D**

D with conventional CVD risk factors

74



**Meta-analysis of D**
Comparison of study weights

D with conventional CVD risk factors

75



**Meta-analysis of C-index change**

C-index change with addition of log CRP

76



**Meta-analysis of C-index change**
Comparison of study weights

C-index change with addition of log CRP

77



**Meta-analysis of D change**

D change with addition of log CRP

78

13

**Meta-analysis of D change**
Comparison of study weights

D change with addition of log CRP

The pooled estimates for the 4 weighting methods respectively are:
0.0414 (p = 0.000), 0.0331 (p = 0.000), 0.0369 (p = 0.000), and 0.0342 (p = 0.000)

79



**Good study-level correlation between absolute values of C-index and D, as well as between changes.**

80



**The large between-study heterogeneity in absolute values of C-index or D is explained by age distribution.**

Size of circle is proportional to 1/variance of estimate

81



**Examination of subgroup effects need necessarily restrict inferences to studies with all levels observed.**

82



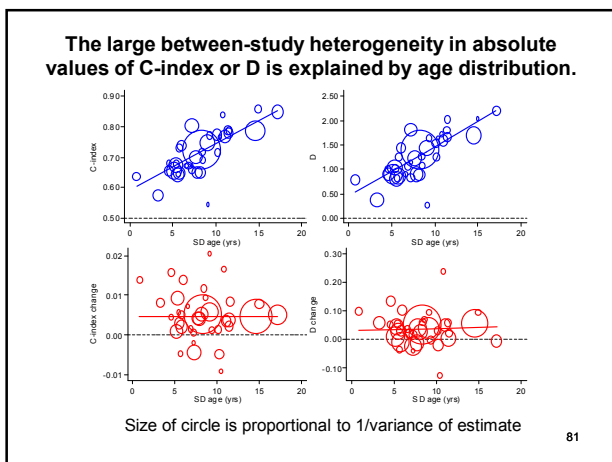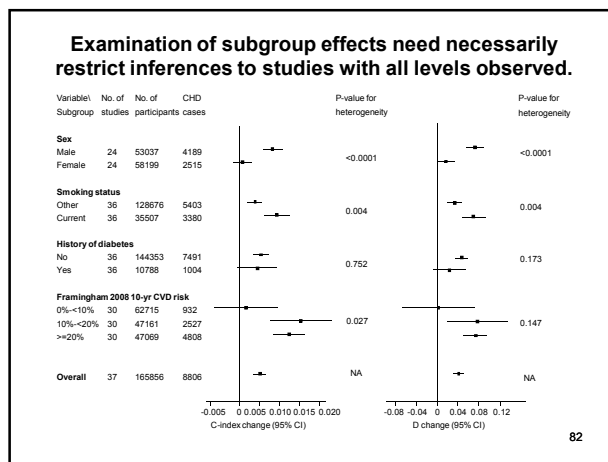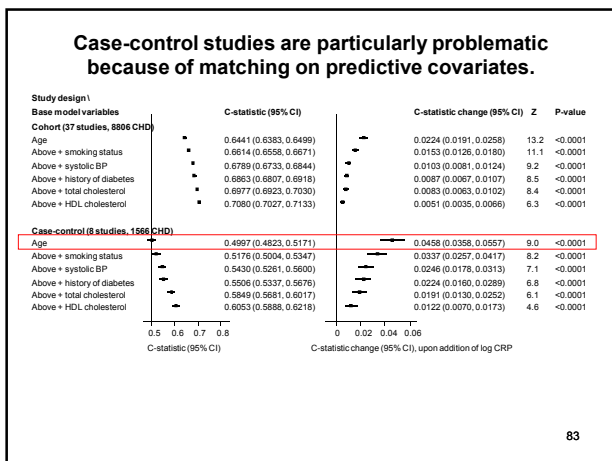**Case-control studies are particularly problematic because of matching on predictive covariates.**

83

# Reclassification

- Methods proposed to give a more clinically relevant assessment of the improvement in risk prediction given by a new marker.

- Examine the net movement of participants between clinically relevant groups of predicted risk (e.g. 10-year CVD risk groups of 0-10%, 10-20%, and ≥20%).

- The Net Reclassification Improvement (NRI)
  - Retrospective NRI (Pencina et al 2008)
  - Prospective NRI
  - Continuous NRI
  
  Pencina et al 2010

84

## 'Retrospective' NRI (Pencina et al 2008)

- Measure of appropriate reclassification (ie: events move up and non-events move down the risk categories)

NRI = P(up|event) - P(down|event) + P(down|non-event) - P(up|non-event)

Net proportion of events appropriately reclassified

Net proportion of non-events appropriately reclassified

85

---

## Reclassification table

Reclassification table for addition of CRP to conventional CVD risk factors (22 studies, 72574 participants):

| Model without CRP | Model with CRP: 10 yr CVD risk (%) | | |
|---|---|---|---|
| 10-yr CVD risk | 0 to <10% | 10% to <20% | ≥20% |
| Observed to have a CVD event within 10 years (n=6848) | | | |
| 0 to <10% | 2424 | 162 | 0 |
| 10% to <20% | 131 | 1684 | 213 |
| ≥20% | 0 | 144 | 2088 |
| Observed to be event free at 10 years (n=65728) | | | |
| 0 to <10% | 53102 | 922 | 0 |
| 10% to <20% | 993 | 6959 | 432 |
| ≥20% | 0 | 403 | 2917 |

86

---

## Summary measures of reclassification

| Addition of CRP to CVD risk factors (22 studies, 72574 participants) | N (%) or Statistic | P-value |
|---|---|---|
| Participants who developed CVD within 10 years | 6846 | |
| Appropriately reclassified | 375 (5.48%) | |
| Inappropriately reclassified | 275 (4.02%) | |
| No change | 6196 (90.51%) | |
| NRI (95% CI), | 1.46% (0.73%, 2.19%) | <0.0001 |
| | | |
| Participants event free at 10 years | 65728 | |
| Appropriately reclassified | 1396 (2.12%) | |
| Inappropriately reclassified | 1354 (2.06%) | |
| No change | 62978 (95.82%) | |
| NRI (95% CI), | 0.06% (-0.09%, 0.22%) | 0.423 |
| | | |
| Overall | | |
| NRI (95% CI) | 1.52% (0.78%, 2.27%) | <0.0001 |
| | | |
| IDI (95% CI) | 0.0036 (0.0028, 0.0043) | <0.0001 |

87

---

## Issue with the 'Retrospective' NRI

- Probabilities estimates are regarding the likelihood of moving up or down through risk categories among events and non-events:

NRI = P(up|event) - P(down|event) + P(down|non-event) - P(up|non-event)

- Need to know whether each participant has or has not had an event by 10 years.

- Cannot include participants whose records are censored before 10 years.

88

---

## Prospective NRI (Pencina et al 2010)

- Pencina et al rearranged the 'retrospective' equation, borrowing from Bayes theorem to give:

$$NRI = \frac{P(events \mid up) \times n_U - P(event \mid down) \times n_D}{n \times P(event)} + \frac{(1 - P(events \mid down)) \times n_D - (1 - P(event \mid up)) \times n_U}{n \times (1 - P(event))}$$

- Now we only need to know whether a person moves up or down a risk category (this is defined for all participants).
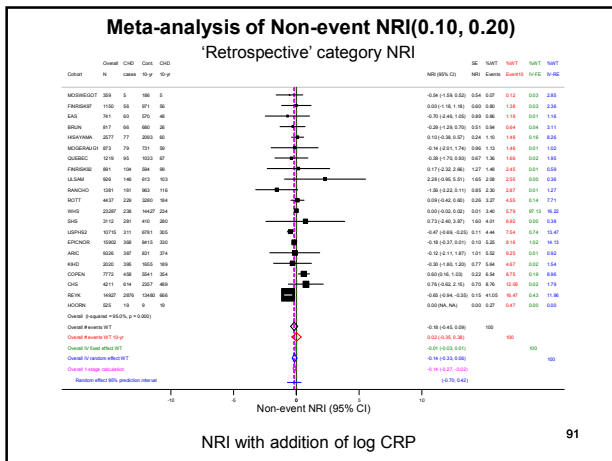
- The quantities P(event | up) and P(event | down) are easily estimated using Kaplan-Meier methods, and hence all censored observations can be included.
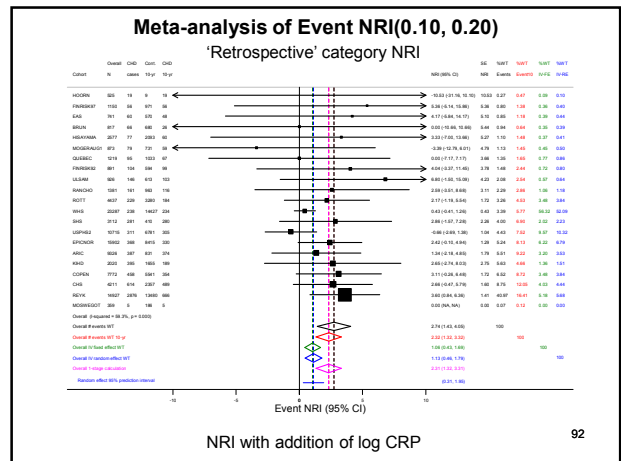
89

---

## ERFC methods: Reclassification

- Reclassification: NRI(0.10, 0.20), NRI(>0), IDI.
  - 10-year risk calculated from stratified Cox model.

  - Category NRI(0.10, 0.20) calculated based on reclassification tables collapsed across studies by observed event status at 10-years (i.e. 1-stage). IDI also calculated by 1-stage method.

  - Sensitivity analysis using 2-stage method (i.e. pooling study-specific NRI(0.10, 0.20) and IDI by meta-analysis).

  - Continuous NRI(>0) calculated within studies based on K-M estimates of 10-year risk among those reclassified up vs. down and then pooled across studies by meta-analysis (i.e. 2-stage).

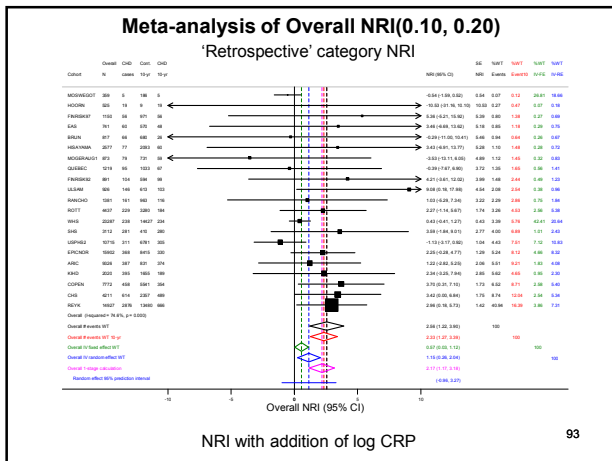  - Appropriate weighting of the study-specific NRI's remains to be determined.

90
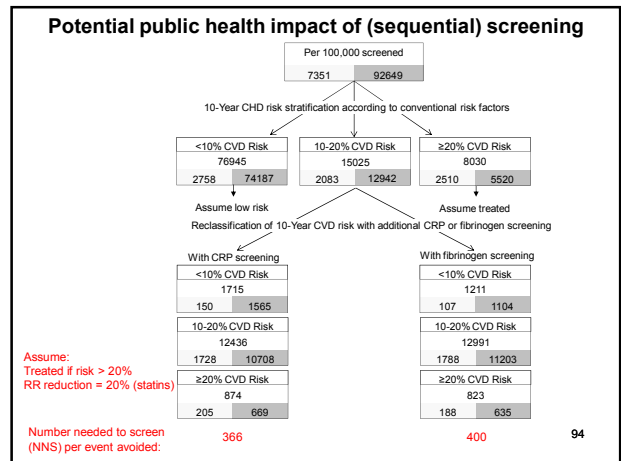
15

## Meta-analysis of Non-event NRI(0.10, 0.20)

'Retrospective' category NRI



Non-event NRI (95% CI)

NRI with addition of log CRP

91

## Meta-analysis of Event NRI(0.10, 0.20)

'Retrospective' category NRI



Event NRI (95% CI)

NRI with addition of log CRP

92

## Meta-analysis of Overall NRI(0.10, 0.20)

'Retrospective' category NRI



Overall NRI (95% CI)

NRI with addition of log CRP

93

## Potential public health impact of (sequential) screening



94

## Remarks

- IPD meta-analysis provides a powerful means for quantifying improvements in risk prediction, but comes with challenges in choice of weighting across studies.

- In particular, uncertainties remain in 2-stage methods for meta-analysis of NRI/IDI because of the separate calculations in Events and Non-Events.
  - What weight is most appropriate for the meta-analysis of the Event and Non-Event NRI's (possibly N and Events)?
  - How should the overall NRI be calculated? Pool the study-specific overall NRI's? or calculate from pooled Event and Non-Event NRI's (possibly the former, but what weights?).
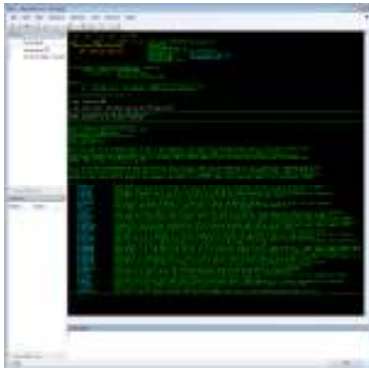  - To what extent should we be concerned of sparse data in small studies?

95

## Stata software for IPD meta-analysis

http://www.phpc.cam.ac.uk/ceu/research/erfc/stata/



96

16

## Installable from within Stata

. net from http://www.phpc.cam.ac.uk/ceu/research/erfc/stata/



97

## Remarks

- Large-scale IPD meta-analysis is helping to powerfully unravel finer details of exposure-disease associations and assessment of improvement in risk prediction than possible in individual studies.

- Methods and software developed in ERFC have been made publicly available.

  – Methods paper: Int J Epidemiol 2010(39)1345-59.

  – Stata software website:
  http://www.phpc.cam.ac.uk/ceu/research/erfc/stata/

98

## Acknowledgements

- Mentors
  – S Thompson, J Danesh

- Statistical collaborators
  – S Thompson, I White, L Pennells, A Wood

- ERFC collaborators
  – http://www.phpc.cam.ac.uk/ceu/research/erfc/studies/

99