

Asian Lung Cancer Absolute Risk Models for lung cancer mortality based on China Kadoorie Biobank

Matthew T. Warkentin, MSc^{1,2}

Martin C. Tammemägi, PhD³

Oswaldo Espin-Garcia, PhD^{2,4}

Sanjeev Budhathoki, PhD¹

Geoffrey Liu, MD^{2,5}

Rayjean J. Hung, PhD^{1,2}

Affiliations:

1. Prosserman Center for Population Health Research, Lunenfeld-Tanenbaum Research Institute, Sinai Health, Toronto, Ontario, Canada
2. Dalla Lana School of Public Health, University of Toronto, Toronto, Ontario, Canada
3. Department of Health Sciences, Brock University, St. Catharines, Ontario, Canada
4. Department of Biostatistics, Princess Margaret Cancer Centre, University Health Network, Toronto, Ontario, Canada
5. Department of Medical Oncology and Hematology, Princess Margaret Cancer Centre, Toronto, Ontario, Canada

Corresponding Author: **Rayjean J. Hung** (rayjean.hung@lunenfeld.ca)

60 Murray St. Toronto, ON M5T 3L9. Canada

Abstract

Background

Lung cancer is the leading cause of cancer mortality globally. Early detection through screening can markedly improve prognosis and prediction models can identify high-risk individuals for screening. However, most models have been developed in North American cohorts of smokers and much less is known about risk factors for never-smokers, which represent a growing proportion of lung cancers, particularly for Asian populations.

Methods

Based on the China Kadoorie Biobank, a population-based prospective cohort study of 512,714 adults age 30-79 recruited between 2004-2008 with up to 12 years of follow-up, we built an Asian Lung Cancer Absolute Risk Model (ALARM) for lung cancer mortality using flexible parametric survival models, separately for ever- and never-smokers. Model performance was evaluated in a 25% held-out test set using the time-dependent area under the receiver operating characteristic curve (AUC) and by comparing the model-predicted and observed risks for agreement (i.e. calibration).

Results

Predictors assessed in the never-smoker lung cancer mortality model were age, sex, household income, lung function, history of emphysema/bronchitis, family history of cancer, personal cancer history, BMI, passive smoking, and indoor air pollution. The ever-smoker model additionally assessed smoking status (former vs. current), duration, intensity, and years since cessation. The 5-year AUC based on the hold-out test set for the never and ever-smoker models were 0.77 (95% CI: 0.73-0.81) and 0.81 (95% CI: 0.79-0.84), respectively. The maximum 5-year risk for never and ever smokers were 2.7% and 14.0%, respectively.

Conclusions

This study is among the first to develop and test risk models specifically for Asian population, separately for never (ALARM-NS) and ever-smokers (ALARM-ES). Our models identify Asian never- and ever-smokers at high-risk of death due to lung cancer with a high degree of accuracy, and may identify those with risks exceeding common eligibility thresholds who would likely benefit from lung cancer screening.

Introduction

Lung cancer is the leading cause of cancer mortality globally. In 2020, there was an estimated 2.2 million incident lung cancers and 1.8 million deaths due to lung cancer, representing 11.4% and 18.0% of all cancer-related incidence and mortality, respectively (1). This corresponds to nearly 1 in 10 cancer diagnoses and 1 in 5 cancer deaths. The five-year survival proportion for lung cancer patients remains poor at only 10-20%, though this varies between countries (1). However, several large randomized trials in populations of predominantly European ancestry have demonstrated a significant reduction in lung cancer mortality when screening with low-dose helical computed tomography (2–5). Identifying high-risk individuals who are likely to benefit from lung cancer screening remains an important public health priority for reducing lung cancer mortality.

The United States Preventive Services Task Force (USPSTF) recommendations were recently updated, and suggest screening adults age 50 to 80 years, with 20 or more pack-year history of smoking, and who have quit smoking no more than 15 years prior (6). Using these criteria would fail to screen any light, long-term former, or never-smokers, which represent a growing proportion of all lung cancer diagnoses, particularly in Asian populations. The proportion of lung cancers among never-smokers varies geographically, with about 15% of lung cancers occurring among never-smokers in North America, and as much as 30-40% in Asian countries (7).

Absolute risk models have been shown to be superior in identifying high-risk individuals for lung cancer screening compared to the USPSTF criteria; however, these models have been primarily developed in North American cohorts of ever-smokers. The widely-used PLCO_{m2012} risk model (8) was developed in a cohort of smokers in the United States, which may not be generalizable to Asian populations due to potential differences in risk factors and baseline risk. Much less is known about the risk factors for never-smokers in Asian populations, within which never-smoker lung cancer is more common than for other racial groups. While approximately 54% of worldwide lung cancers occurred in Asian countries (1), currently there is no validated risk prediction model specifically for Asian populations.

The goal of the current study was to develop and evaluate a lung cancer mortality risk prediction model, specifically for Asian populations, with separate models developed for ever- and never-smokers.

Methods

Study Participants

This study used data from the China Kadoorie Biobank (CKB) which has been described in detail previously (9). In brief, the CKB is a population-based prospective cohort of 512,714 adults age 30 to 79 recruited between 2004 and 2008 from 10 geographically defined regions in China with 5.1 million person-years of follow-up and detailed collections of epidemiologic data. We excluded any participants with a personal history of lung cancer within 5 years of baseline. Mortality status was collected through electronic linkage to regional mortality registries with up to 12 years of follow-up. Deaths due to lung cancer was the primary outcome in this study.

Statistical Analysis

Absolute risk models

To estimate the probability of lung cancer mortality (deaths attributed to lung cancer), we modeled the cause-specific hazards for death from lung cancer. We used flexible parametric survival models (i.e., Royston-Parmar models (10)) on the cumulative hazard scale to estimate baseline hazards and hazard ratios for predictor effects (e.g., clinico-epidemiological and spirometry data), using complete-case analysis separately for ever- and never-smokers. Models were fit using the `flexsurv` package in R (11).

Within these absolute risk models, time-since-entry was used as the timescale using restricted (natural) cubic splines with internal knots placed at the quantiles of the log uncensored cause-specific event-time distributions. We estimate the 5-year probability of lung cancer mortality conditional on a set of risk factors (See Supplementary Methods for more details).

Asian Lung cancer Absolute Risk Models (ALARM) were built separately for never (ALARM-NS) and ever-smokers (ALARM-ES). We used stratified random sampling to split the data into training and testing sets, maintaining similar proportions of the outcome, separately for never- and ever-smokers. Models were fitted in the training (75%) data, and the held-out testing (25%) data were used to estimate out-of-sample performance for internal model validation.

For the lung cancer mortality models, potential predictors were selected based on their known associations with lung cancer or by improving model performance. ALARM-NS included age, sex, body mass index, family history of cancer, personal cancer history (excluding lung cancers within 5 years of enrollment), lung function (forced expiratory volume in 1 second (FEV1) / forced vital capacity (FVC)), personal history of emphysema or bronchitis, household income (5 levels), passive smoking, and cooking fuel exposure. ALARM-ES included these variables with the addition of smoking status (current vs. former), smoking duration (years), smoking intensity (cigarettes / day), and time since quitting (years). Cooking fuel exposure was computed as a composite of cooking fuel type (gas or electricity vs. all others) and duration of usage (years of usage at current and previous residences). Numeric variables were evaluated for potential non-linear relationships and time-dependent effects. Age was included in the final models as the logarithm of age, and BMI was included as quadratic (BMI and BMI²). We report the hazard ratios (HR) and 95% confidence intervals (CI). All statistical analyses were done using R version 4.0.5 (12).

Model performance

The performance of the prediction models were assessed in two complementary ways: (1) how closely model-predicted risks corresponded to observed risks (calibration) and, (2) the ability of the model to assign higher predicted risks to those who died of lung cancer, or died of lung cancer at an earlier time (discrimination).

Model calibration was assessed by comparing the observed five-year risks to model-predicted (expected) five-year risks, separately for never-smokers and ever-smokers. Graphical assessment was done by comparing risks within binned risk groups. The observed five-year risks were estimated based on the complement of the non-parametric Kaplan-Meier (KM) estimator for survival (13). The model-predicted five-year risks were estimated based on the absolute risk models described above. Ideal calibration corresponds to points falling along the diagonal line in a calibration plot with a slope of 1 (i.e. the identity line). The ratio of expected to observed deaths (E/O) and difference in expected and observed deaths (E-O) per 100,000 are reported for the never and ever-smoker models.

Discrimination was measured by the area under the time-dependent receiver operating characteristic curve (AUC). We used the definition of the time-dependent ROC described in Blanche *et al* (15). We

constructed 95% confidence intervals for the time-dependent AUC and calibration metrics using a percentile-based approach with 500 bootstrap resamples. Calibration and discrimination metrics are reported for the hold-out test data only.

Results

Basic demographic characteristics for the CKB cohort according to vital status based on the total length of follow-up are presented in **Table 1**. In general, those who died from lung cancer were older at baseline (62.2 vs. 52.0 years), were more likely to be male (64% vs. 41%), and were more likely to be current or former smokers with extensive smoking history, when compared with those still alive.

In total, two separate parametric survival models were fitted to form the absolute risk models for never (ALARM-NS) and ever-smokers (ALARM-ES). The hazard ratios and 95% confidence intervals for the lung cancer mortality models are reported in **Table 2**.

Age (on the log scale) increased the risk of lung cancer mortality for both never and ever smokers. BMI was included in the model as quadratic terms to capture the U-shaped relationship between BMI and lung cancer mortality. Female sex was found to be protective for lung-cancer mortality in never-smokers (HR: 0.84, 95% CI: 0.70-1.00), but a modest risk factor in ever-smokers (HR: 1.18, 95% CI: 0.96-1.44). Family history of any type of cancer was a risk factor for ever-smoker death due to lung cancer (HR: 1.29, 95% CI: 1.12-1.48), but not for never-smokers. Personal cancer history increased hazard of lung cancer mortality in never (HR: 1.30, 95% CI: 0.62-2.74) and ever smokers (HR: 1.95, 95% CI: 1.15-3.31). For every 5% increase in lung function performance (FEV1/FVC), the hazard of lung cancer mortality was reduced by 6% and 2%, for never and ever smokers, respectively. In addition, a self-reported history of emphysema or bronchitis had HR of 1.03 (95% CI: 0.78-1.37) and 1.30 (95% CI: 1.08-1.57) for never and ever-smokers, respectively. Cooking fuel exposure was not found to improve model performance and was therefore not included from the final models.

ALARM-ES further included several smoking variables. When compared to current smokers, former smokers had a lower hazard of lung cancer mortality (HR: 0.85, 95% CI: 0.71-1.01). For every 5 additional years of smoking, the hazard of lung cancer mortality were increased by 19% (HR: 1.19, 95% CI: 1.15-

1.24). For every additional 10 cigarettes smoked per day (approximately equivalent to half a pack), the hazard of lung cancer mortality was increased by 20% (HR: 1.20, 95% CI: 1.15-1.26). Non-linear relationships and two-way interactions were explored for smoking variables but did not contribute to model improvement and were not included in the final models.

We estimated the five-year absolute risk of lung cancer mortality (i.e., cumulative incidence of death from lung cancer), conditional on a participant's risk factor profile. The distribution of risks for never- and ever-smokers are presented in **Supplemental Figure 1**. The maximum five-year predicted risk of lung cancer mortality was 14.0% for ever-smokers and 2.7% for never-smokers. According to our models, 9.2% of ever-smokers and <1% of never-smokers in the CKB would be eligible for screening assuming a 1.5% five-year risk threshold.

Calibration plots comparing model-predicted and observed five-year risks across risk quantiles, separately for ALARM-NS and ALARM-ES, are presented in **Supplemental Figure 2**. Both models show very good calibration in the hold-out test data. Additional calibration metrics overall and by key subgroups are presented in Supplementary Table 3. The five-year time-dependent AUC in the hold-out test set 0.77 (95% CI: 0.73-0.81) and was 0.81 (95% CI: 0.79-0.84) for ALARM-NS and ALARM-ES models, respectively. The time-dependent receiver operating characteristic curves for the 25% held-out test data are presented in **Figure 1**.

Absolute risk trajectories for average and high risk current and former smokers are presented in **Figure 2** according to lung function performance and smoking intensity. Lung cancer mortality risks are higher for those with suboptimal lung function and higher smoking intensity. Absolute risk trajectories for never-smokers are presented in **Figure 3**, separately for men and women according to lung function. Lung cancer mortality risk is higher for suboptimal lung function and for men. Contour plots for absolute risk according to smoking intensity, age at smoking initiation, and lung function for a theoretical average or high risk current and former smokers are shown in **Supplemental Figure 3**. Five-year absolute risks between 13.5% and 15.0% are observed for the highest risk profile. Contour plots for a theoretical average or high risk never-smoker according to age, FEV1/FVC, and sex are shown in **Figure 4**, with risks as high as 2.75% to 3.00% for the highest risk profile.

We compared our absolute risk models to the recently-updated USPSTF criteria for lung cancer screening (see **Table 3**). Using the USPSTF criteria, 35.5% of CKB ever-smokers would be eligible for lung cancer screening. To compare our model to the USPSTF criteria, we applied the five-year risk threshold that produced an equivalent specificity to the USPSTF criteria - our model (ALARM-ES) demonstrated an improved sensitivity (82.0% vs. 68.6% based on USPSTF), positive predictive value (1.2% vs. 1.0%), and negative predictive value (99.9% vs. 99.7%), while selecting a similar proportion of the population for screening. At a five-year risk threshold that has an equivalent sensitivity to the USPSTF criteria, ALARM-ES would identify 12% fewer ever-smokers for screening (23.9% vs. 35.5% based on USPSTF).

We applied the established Lung Cancer Death Risk Assessment Tool (LCDRAT) ([16](#)) to the CKB ever-smokers to assess how our model (ALARM-ES) performed when compared against a validated ever-smoker lung cancer mortality model. Details of how LCDRAT was applied and evaluated in the CKB are described in the Supplemental Methods. Calibration and discrimination statistics are presented in **Supplementary Table 3**. LCDRAT and LCDRAT Constrained had similar discriminative performance (i.e. AUC) as ALARM-ES, however, ALARM-ES had superior calibration. The expected/observed death ratios were 1.07 (95% CI: 0.96-1.22) in ALARM-ES and 1.57 (95% CI: 1.46-1.69) in LCDRAT, and the differences between expected and observed deaths (per 100,000) were 37.0 (-23.33-101.46) and 288.02 (95% CI: 250.10-326.47), respectively. This is not surprising as risk models developed in North American cohorts may not be well calibrated in Asian populations, which are reflected in the calibration statistics, despite good overall discrimination.

Discussion

We established risk prediction models based on a relatively simple set of predictors that can be ascertained during a routine physician visit, which can accurately discriminate high and low risk individuals for lung cancer mortality in an Asian population better than the US-based current recommended screening guidelines. These are some of the first prediction models for lung cancer mortality that were developed and evaluated specifically for an Asian population and separately for never and ever-smokers. Our model has shown superior calibration compared to the established lung cancer mortality model (LCDRAT), while maintaining a comparable model accuracy. At a risk threshold that matches the specificity of the USPSTF criteria, ALARM-ES has better sensitivity, PPV, and NPV than USPSTF criteria, while selecting a similar proportion of the population for lung cancer screening.

Despite the global reductions in smoking prevalence in most parts of the world (17), lung cancers among never-smokers is an increasingly important public health concern (18). The proportion of lung cancers in never smokers has been increasing and lung cancer in never smoker is one of the most common cancers (19). It is estimated that the epicenter of lung cancer in the next few decades would be within Asian countries, where a substantial proportion of lung cancers would occur in never smokers. This highlights the importance of having risk models developed specifically for this population. The International Association of Lung Cancer Study (IASLC) recently released a statement specifically focused on never smokers, in which it emphasizes the importance of risk-based screening for never smokers (18). While widespread screening of never-smokers at the current stage would not be effective, the development and validation of risk predictions will play a critical role (18). The report encourages and recommends modelling studies focused on lung cancer risk estimation for never-smokers in order to eventually realize the benefit of risk-based screening in this group (18).

Furthermore, it is well known that lung cancers occurring in never smokers represent a distinct form of the disease (19,20). Lung cancers occurring among never-smoker are more commonly adenocarcinomas, with higher proportions of EGFR mutations, which makes them highly targetable by therapeutics leading to improved prognosis (21); although these therapeutic agents are typically administered when the disease is at a later, incurable stage, the majority of lung cancer patients with EGFR-mutations will eventually become incurable at some point in their cancer course. This contributes, in part, to the

improved survival observed among never-smokers (21). Based on recent findings from the ADAURA trial, there should be an even stronger push to identify these cancers at an earlier resectable stage, where the future use of EGFR tyrosine kinase inhibitors can improve disease free survival (22). In the absence of primary smoking history as a risk factor, statistical models based on a constellation of risk factors will be required in order to identify high-risk never-smokers for screening. The ability to identify and screen high-risk never-smokers for lung cancer is an important public health concern.

In this study, we developed and evaluated risk models for estimating the five-year absolute risk of lung cancer mortality in the held-out test set. Lung cancers have previously been found to affect never-smoking women disproportionately more than men (21,23), though women are still observed to have a survival advantage (24), which is consistent with our findings. Due to limitations in the data, we were unable to determine a participant's family history of lung cancer, so we used family history of any cancer, which is expected to be a relatively weak proxy. Exposure to cooking fuel has previously been identified as a risk factor for lung cancer (25,26), particularly among Asian women (27,28). However, in our study we did not observe an association between cooking fuel exposure and lung cancer mortality. We believe this may be due, in part, to imprecise measurements of this putative risk factor. Similar to passive tobacco exposure (i.e. second-hand smoke), cooking fuel exposure is subject to recall bias and thus it is difficult to accurately measure cumulative lifetime exposure.

We identified four previous studies which developed or adapted risk models for never-smokers, and two were based in Asian populations. The PLCOall2014 model was developed in both ever and never-smoker population, and was adapted for use in never-smokers by removing smoking-related predictors (29). We applied the PLCOall2014 to the CKB, and our model achieved higher AUC for never smokers (0.77 versus 0.76 for PLCOall2014) and ever-smokers (0.81 vs. 0.75 based on PLCOall2014). Details of this comparison are described in the Supplementary Methods. A second study was developed in a predominantly Caucasian cohort (UK Biobank) (30). A third study was developed in a Taiwanese cohort and models were built separately for never, light, and heavy smokers (31). This study achieved a training-set AUC of 0.806 (95% CI: 0.790-0.819) among never-smokers, but included several serological protein biomarkers unavailable in the CKB that are not routinely tested in health populations and therefore could not be validated. The authors presented limited evidence of model calibration across risk groups and by smoking status. A fourth model was developed for never-smoker females only in an Asian case-control

study and achieved an AUC of 0.71 (95% CI: 0.66-0.77) (32). However, this study was not based on prospective data, and would require the collection of genetic data (i.e., genotyping) which may be prohibitive for widespread application. The model without the 9 genetic variants performed moderately worse (AUC=0.69, 95% CI: 0.64-0.75).

Based on ALARM-NS, fewer than 1% of the never-smokers in CKB reached a five-year risk threshold of at least 1.5% to be eligible lung cancer screening, with the highest-risk never smoker achieving a risk of 2.7%. The marginal age-specific lung cancer mortality rates for men and women were lower in the CKB than those observed in the general Chinese population, based on Global Burden of Disease 2019 mortality estimates (see Supplemental Figure 4). This may be due, in part, to a “healthy volunteer” effect. This phenomenon occurs when those persons who volunteer for a study are healthier and not fully representative of the general population. As such, we anticipate the actual risk distribution in the general Chinese population would be higher than what is observed in CKB (Supplementary Figure 1). It is possible that potential measurement error of important predictors (e.g. indoor air pollution) has led to an underestimate of absolute risks. A larger proportion of individuals would be at a higher risk for lung cancer mortality and may be eligible for screening based on applying our model to the general Chinese population.

In summary, we developed absolute risk models that accurately identify both ever- and never-smokers at high-risk for lung cancer-related mortality. These models were developed exclusively using prospective data collected in a large Asian population, and models were fitted separately based on smoking history (ALARM-NS and ALARM-ES). Our models discriminated high and low risk for lung cancer death similarly well to an established lung cancer mortality model (i.e. LCDRAT), but demonstrated much better absolute risk calibration for an Asian population. In the future, these models may be useful for identifying a subset of the Asian population at high-risk for lung cancer mortality who may benefit from lung cancer screening. The next step will be to externally validate our models in an independent cohort.

Funding

This work was supported by Canadian Institutes of Health Research (FDN 167273) and the National Institutes of Health (U19 CA203654).

Notes

Role of the funders

The funding sources had no role in the conception, design, implementation, analysis, or interpretation of the study, or writing of the report, or the decision to submit the report for publication.

Disclosures

The authors report no potential conflicts of interests.

Contributors

RJH conceived the study and obtained the funding. MTW and RJH contributed equally to the design of the study and co-led the writing of the article. MTW led the statistical analysis with scientific input from RJH, OEG, and MCT. All coauthors contributed to the result interpretations, critical review of the manuscript, and approval of the final version.

References

1. Sung H, Ferlay J, Siegel RL, et al. Global cancer statistics 2020: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA: a cancer journal for clinicians*. Published online 2021.
2. Koning HJ de, Aalst CM van der, Jong PA de, et al. Reduced lung-cancer mortality with volume CT screening in a randomized trial. *New England journal of medicine*. 2020;382(6):503-513.
3. National Lung Screening Trial Research Team. Reduced lung-cancer mortality with low-dose computed tomographic screening. *New England Journal of Medicine*. 2011;365(5):395-409.
4. National Lung Screening Trial Research Team. Lung cancer incidence and mortality with extended follow-up in the national lung screening trial. *Journal of Thoracic Oncology*. 2019;14(10):1732-1742.
5. Pastorino U, Silva M, Sestini S, et al. Prolonged lung cancer screening reduced 10-year mortality in the MILD trial: New confirmation of lung cancer screening efficacy. *Annals of Oncology*. 2019;30(7):1162-1169.
6. Moyer VA. Screening for lung cancer: US preventive services task force recommendation statement. *Annals of internal medicine*. 2014;160(5):330-338.
7. Toh CK, Lim WT. Lung cancer in never-smokers. *Journal of clinical pathology*. 2007;60(4):337-340.
8. Tammemägi MC, Katki HA, Hocking WG, et al. Selection criteria for lung-cancer screening. *New England Journal of Medicine*. 2013;368(8):728-736.
9. Chen Z, Chen J, Collins R, et al. China kadoorie biobank of 0.5 million people: Survey methods, baseline characteristics and long-term follow-up. *International journal of epidemiology*. 2011;40(6):1652-1666.
10. Royston P, Parmar MK. Flexible parametric proportional-hazards and proportional-odds models for censored survival data, with application to prognostic modelling and estimation of treatment effects. *Statistics in medicine*. 2002;21(15):2175-2197.
11. Jackson C. flexsurv: A platform for parametric survival modeling in R. *Journal of Statistical Software*. 2016;70(8):1-33. doi:[10.18637/jss.v070.i08](https://doi.org/10.18637/jss.v070.i08)
12. R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing; 2021. <https://www.R-project.org/>
13. Gerds TA, Ozenne B. *riskRegression: Risk Regression Models and Prediction Scores for Survival Analysis with Competing Risks.*; 2020. <https://CRAN.R-project.org/package=riskRegression>

14. Blanche P, Dartigues JF, Jacqmin-Gadda H. Estimating and comparing time-dependent areas under receiver operating characteristic curves for censored event times with competing risks. *Statistics in medicine*. 2013;32(30):5381-5397.
15. Blanche P, Dartigues JF, Jacqmin-Gadda H. Estimating and comparing time-dependent areas under receiver operating characteristic curves for censored event times with competing risks. *Statistics in Medicine*. 2013;32(30):5381-5397. <http://onlinelibrary.wiley.com/doi/10.1002/sim.5958/full>
16. Katki HA, Kovalchik SA, Berg CD, Cheung LC, Chaturvedi AK. Development and validation of risk models to select ever-smokers for CT lung cancer screening. *Jama*. 2016;315(21):2300-2311.
17. Geneva: World Health Organization. *WHO Global Report on Trends in Prevalence of Tobacco Smoking 2000-2025, Second Edition.*; 2018.
18. Kerpel-Fronius A, Tammemägi M, Cavic M, et al. Screening for lung cancer in individuals who never smoked: An international association for the study of lung cancer early detection and screening committee report. *Journal of Thoracic Oncology*. 2022;17(1):56-66.
19. Sun S, Schiller JH, Gazdar AF. Lung cancer in never smokers—a different disease. *Nature Reviews Cancer*. 2007;7(10):778-790.
20. Couraud S, Zalcman G, Milleron B, Morin F, Souquet PJ. Lung cancer in never smokers—a review. *European journal of cancer*. 2012;48(9):1299-1311.
21. Subramanian J, Govindan R. Lung cancer in never smokers: A review. *Journal of clinical oncology*. 2007;25(5):561-570.
22. Wu YL, Herbst RS, Mann H, Rukazenzov Y, Marotti M, Tsuboi M. ADAURA: Phase III, double-blind, randomized study of osimertinib versus placebo in EGFR mutation-positive early-stage NSCLC after complete surgical resection. *Clinical lung cancer*. 2018;19(4):e533-e536.
23. Freedman ND, Leitzmann MF, Hollenbeck AR, Schatzkin A, Abnet CC. Cigarette smoking and subsequent risk of lung cancer in men and women: Analysis of a prospective cohort study. *The lancet oncology*. 2008;9(7):649-656.
24. Thun MJ, Hannan LM, Adams-Campbell LL, et al. Lung cancer occurrence in never-smokers: An analysis of 13 cohorts and 22 cancer registry studies. *PLoS Med*. 2008;5(9):e185.
25. Kurmi OP, Arya PH, Lam KBH, Sorahan T, Ayres JG. Lung cancer risk and solid fuel smoke exposure: A systematic review and meta-analysis. *European Respiratory Journal*. 2012;40(5):1228-1237.
26. Lissowska J, Bardin-Mikolajczak A, Fletcher T, et al. Lung cancer and indoor pollution from heating and cooking with solid fuels: The IARC international multicentre case-control study in eastern/central europe and the united kingdom. *American journal of epidemiology*. 2005;162(4):326-333.

27. Xue Y, Jiang Y, Jin S, Li Y. Association between cooking oil fume exposure and lung cancer among chinese nonsmoking women: A meta-analysis. *OncoTargets and therapy*. 2016;9:2987.
28. Chen TY, Fang YH, Chen HL, et al. Impact of cooking oil fume exposure and fume extractor use on lung cancer risk in non-smoking han chinese women. *Scientific reports*. 2020;10(1):1-10.
29. Tammemaegi MC, Church TR, Hocking WG, et al. Evaluation of the lung cancer risks at which to screen ever-and never-smokers: Screening rules applied to the PLCO and NLST cohorts. *PLoS Med*. 2014;11(12):e1001764.
30. Muller DC, Johansson M, Brennan P. Lung cancer risk prediction model incorporating lung function: Development and validation in the UK biobank prospective cohort study. In: American Society of Clinical Oncology; 2017.
31. Wu X, Wen CP, Ye Y, et al. Personalized risk assessment in never, light, and heavy smokers in a prospective cohort in taiwan. *Scientific reports*. 2016;6(1):1-9.
32. Chien LH, Chen CH, Chen TY, et al. Predicting lung cancer occurrence in never-smoking females in asia: TNSF-SQ, a prediction model. *Cancer Epidemiology and Prevention Biomarkers*. 2020;29(2):452-459.

Table 1. Distribution of demographic characteristics of China Kadoorie Biobank based on mortality status. Means and standard deviations are reported for numeric variables, and frequencies and proportions for categorical variables.

	Alive No. (%)	Lung cancer mortality No. (%)	Total No. (%)
N	509,885	2,754	512,639
Age (years)			
mean [sd]	52.0 [10.7]	62.2 [9.1]	52 [10.7]
Length of follow-up (years)			
mean [sd]	10.0 (1.8)	4.9 [2.3]	9.9 [1.8]
Sex			
Female	301,506 (59%)	991 (36%)	302,497 (59%)
Male	208,379 (41%)	1,763 (64%)	210,142 (41%)
Family history of cancer			
No	423,079 (83%)	2,253 (82%)	425,332 (83%)
Yes	86,806 (17%)	501 (18%)	87,307 (17%)
Personal cancer history			
No	507,412 (>99%)	2,724 (99%)	510,136 (>99%)
Yes	2,473 (<1%)	30 (1%)	2,503 (<1%)
Emphysema or bronchitis			
No	496,798 (97%)	2,553 (93%)	499,351 (97%)
Yes	13,087 (3%)	201 (7%)	13,288 (3%)
Smoking status			
Never	345,557 (68%)	1,060 (38%)	346,617 (68%)
Former	37,553 (7%)	415 (15%)	37,968 (7%)
Current	126,775 (25%)	1,279 (46%)	128,054 (25%)
Smoking intensity (cigs/day) ^a			
mean [sd]	17.8 [10.8]	19.0 [11.2]	17.8 [10.8]
Smoking duration (years) ^a			
mean [sd]	28.4 [11.6]	38.8 [11.1]	28.5 [11.6]
Years since cessation (years) ^b			
mean [sd]	8.7 [8.5]	8.2 [8.5]	8.7 [8.5]
FEV1/FVC (%)			
mean [sd]	84.5 [8.5]	80.6 [10.5]	84.5 [8.5]
Body mass index (kg/m ²)			
mean [sd]	23.7 [3.4]	22.8 [3.6]	23.7 [3.4]
Household income ^c			
<2500 yuan	15,401 (3%)	136 (5%)	15,537 (3%)
2,500-4,999 yuan	34,388 (7%)	236 (9%)	34,624 (7%)
5,000-9,999 yuan	94,059 (18%)	502 (18%)	94,561 (18%)
10,000-19,999 yuan	148,103 (29%)	829 (30%)	148,932 (29%)
20,000-34,999 yuan	126,034 (25%)	647 (23%)	126,681 (25%)
≥35,000 yuan	91,900 (18%)	404 (15%)	92,304 (18%)
Cooking fuel exposure ^d			
None	203,134 (40%)	1,242 (45%)	204,376 (40%)
Low	80,385 (16%)	443 (16%)	80,828 (16%)
Medium	149,649 (30%)	602 (22%)	150,251 (29%)
High	76,717 (14%)	467 (17%)	77,184 (15%)

Abbreviations: FEV1 = forced-expiratory volume, 1-second; FVC = forced vital capacity; sd = standard deviation.

^a Average smoking intensity and number of years smoking (duration) includes current and former smokers; never smokers are excluded.

^b Years since smoking cessation only includes former smokers; current and never smokers are excluded.

^c Approximate equivalents in USD (\$), rounded to the nearest dollar: <384, 384-767, 767-1,537, 1,537-3,075, 3,075-5,381, and ≥5381.

^d Cooking fuel exposure groups were formed based on 25th and 75th percentiles of the cumulative cooking exposure distribution. Cumulative cooking exposure was calculated based on the self-reported frequency of exposure to potentially harmful cooking fuels (i.e., coal, wood, or other, as compared to gas or electricity).

Table 2. Estimates from lung cancer mortality flexible parametric survival models fit separately for never-smokers (ALARM-NS) and ever-smokers (ALARM-ES). Lung cancer cause-specific hazard ratios with 95% confidence intervals and beta coefficients (log hazard ratios) with standard errors are reported.

	ALARM-NS		ALARM-ES	
	HR (95% CI)	Beta (SE)	HR (95% CI)	Beta (SE)
Age at entry (natural log) ^a	—	4.7174 (0.2200)	—	4.1154 (0.2936)
Sex (Female vs. Male)	0.84 (0.70-1.00)	-0.1744 (0.0902)	1.18 (0.96-1.44)	0.1615 (0.1039)
Family history of cancer (Yes vs. No)	0.88 (0.73-1.07)	-0.1245 (0.1000)	1.29 (1.12-1.48)	0.2539 (0.0713)
Personal history of cancer (Yes vs. No)	1.30 (0.62-2.74)	0.2643 (0.3803)	1.95 (1.15-3.31)	0.6679 (0.2706)
FEV1/FVC (per 5%)	0.94 (0.90-0.99)	-0.0583 (0.0253)	0.98 (0.95-1.02)	-0.0167 (0.0186)
Personal history of emphysema or bronchitis (Yes vs. No)	1.03 (0.78-1.37)	0.0314 (0.1428)	1.30 (1.08-1.57)	0.2646 (0.0967)
Household income (per level)	0.98 (0.93-1.03)	-0.0174 (0.0263)	1.08 (1.03-1.12)	0.0735 (0.0210)
Body mass index (kg/m ²) ^b				
BMI	—	-0.1976 (0.1088)	—	-0.1893 (0.1550)
BMI ²	—	0.0036 (0.0020)	—	0.0032 (0.0030)
Smoking status (Former vs. Current)	—	—	0.85 (0.71-1.01)	-0.1673 (0.0897)
Smoking duration (per 5 years)	—	—	1.19 (1.15-1.24)	0.1761 (0.0201)
Smoking intensity (per 10 cigarettes / day)	—	—	1.20 (1.15-1.26)	0.1822 (0.0233)
Years since cessation (per 5 years)	—	—	1.04 (0.97-1.13)	0.0436 (0.0404)

Abbreviations: ALARM = Asian Lung cancer Absolute Risk Model; CI = confidence interval; ES = Ever smoker; FEV1 = forced-expiratory volume, 1-second; FVC = forced vital capacity; HR = hazard ratio; NS = never smoker.

^a Age (in years) was modeled on the natural-log scale. We have excluded the hazard ratio as it is difficult to interpret without inverting the transformation.

^b Body mass index (in kg/m²) was modeled as a quadratic relationship (i.e. BMI and BMI²). We have excluded the hazard ratios as it is difficult to interpret without considering both effects jointly.

Table 3. Comparison of ALARM-ES and LCDRAT against the United States Preventive Services Task Force (USPSTF) 2021 criteria for lung cancer screening. The USPSTF recommends annual screening for adults age 50-80 years, with 20 pack-year smoking history, and smoking cessation less than 15 years prior for former smokers. We compared the ever-smoker models two ways: (1) using a risk threshold that matches the specificity of USPSTF criteria, and (2) using a risk threshold that matches the sensitivity of USPSTF criteria.

	Specificity	Sensitivity	PPV	NPV	% eligible for screening
USPSTF-2021	65.5%	68.6%	1.0%	99.7%	35.5%
ALARM-ES					
Matched on specificity	Same as USPSTF	82.0%	1.2%	99.9%	36.4%
Matched on sensitivity	78.0%	Same as USPSTF	1.6%	99.8%	23.9%
LCDRAT					
Matched on specificity	Same as USPSTF	82.2%	1.1%	97.7%	36.5%
Matched on sensitivity	78.0%	Same as USPSTF	1.4%	97.0%	23.8%

Abbreviations: ALARM = Asian Lung cancer Absolute Risk Model; LCDRAT = Lung Cancer Death Risk Assessment Tool; NPV = Negative predictive value; PPV = Positive predictive value; USPSTF = United States Preventive Services Task Force, USPSTF.

Note: Time-dependent sensitivity, specificity, PPV, and NPV for USPSTF-2021 and ALARM-ES were based on inverse-probability of censoring weighted (IPCW) estimates for a cumulative/dynamic definition of cases and control proposed by Blanche *et al.* (2013). For LCDRAT, the IPCW estimates were based on the time-dependent Sens/Spec/PPV/NPV proposed by Blanche *et al.* (2013) which extends the cumulative/dynamic definition of cases and control for the competing risks setting.

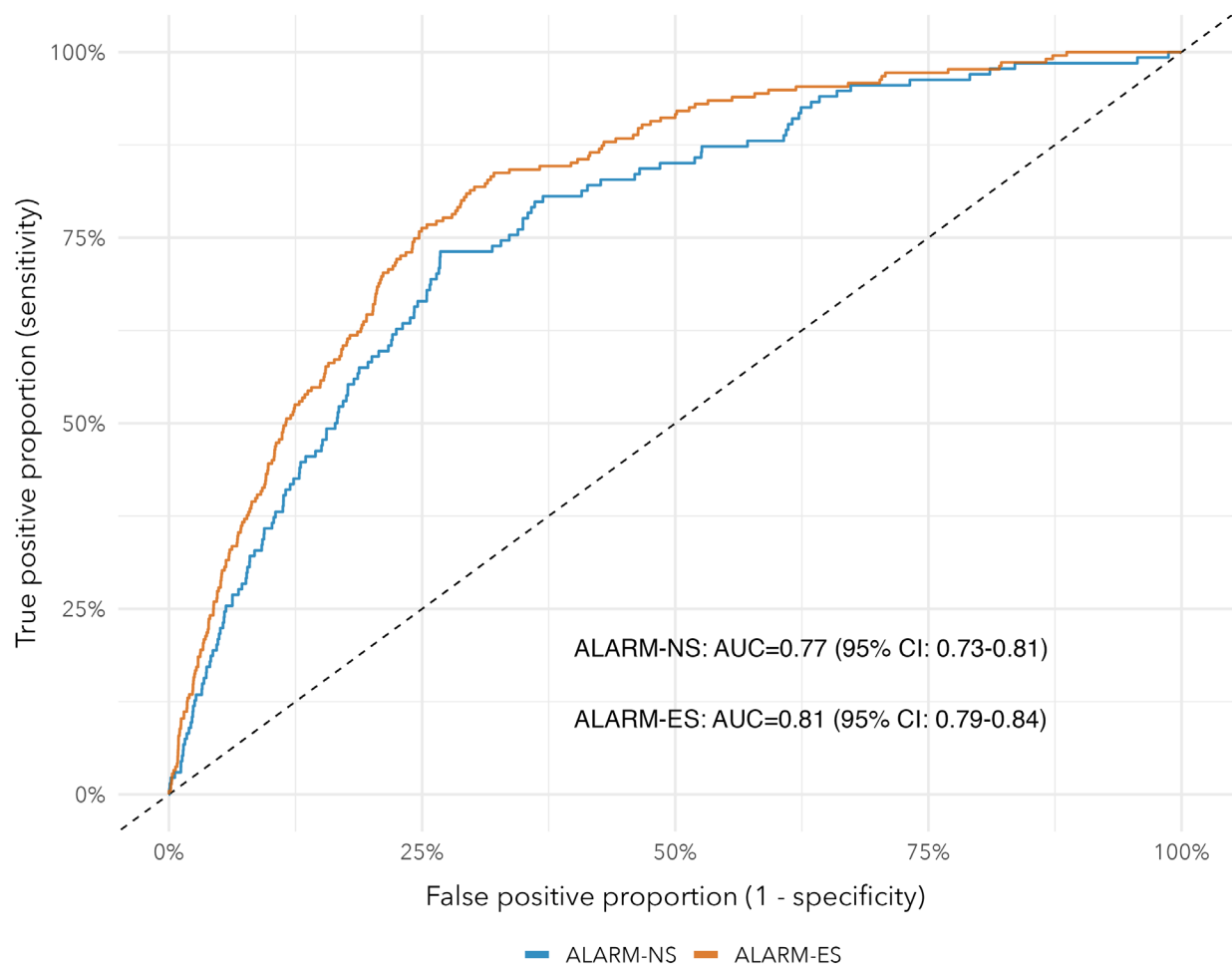


Figure 1: Five-year receiver operating characteristic (ROC) curves and area under the curves (AUC) for the 30% hold-out test data, separately for ALARM-NS (never-smokers) and ALARM-ES (ever-smokers).

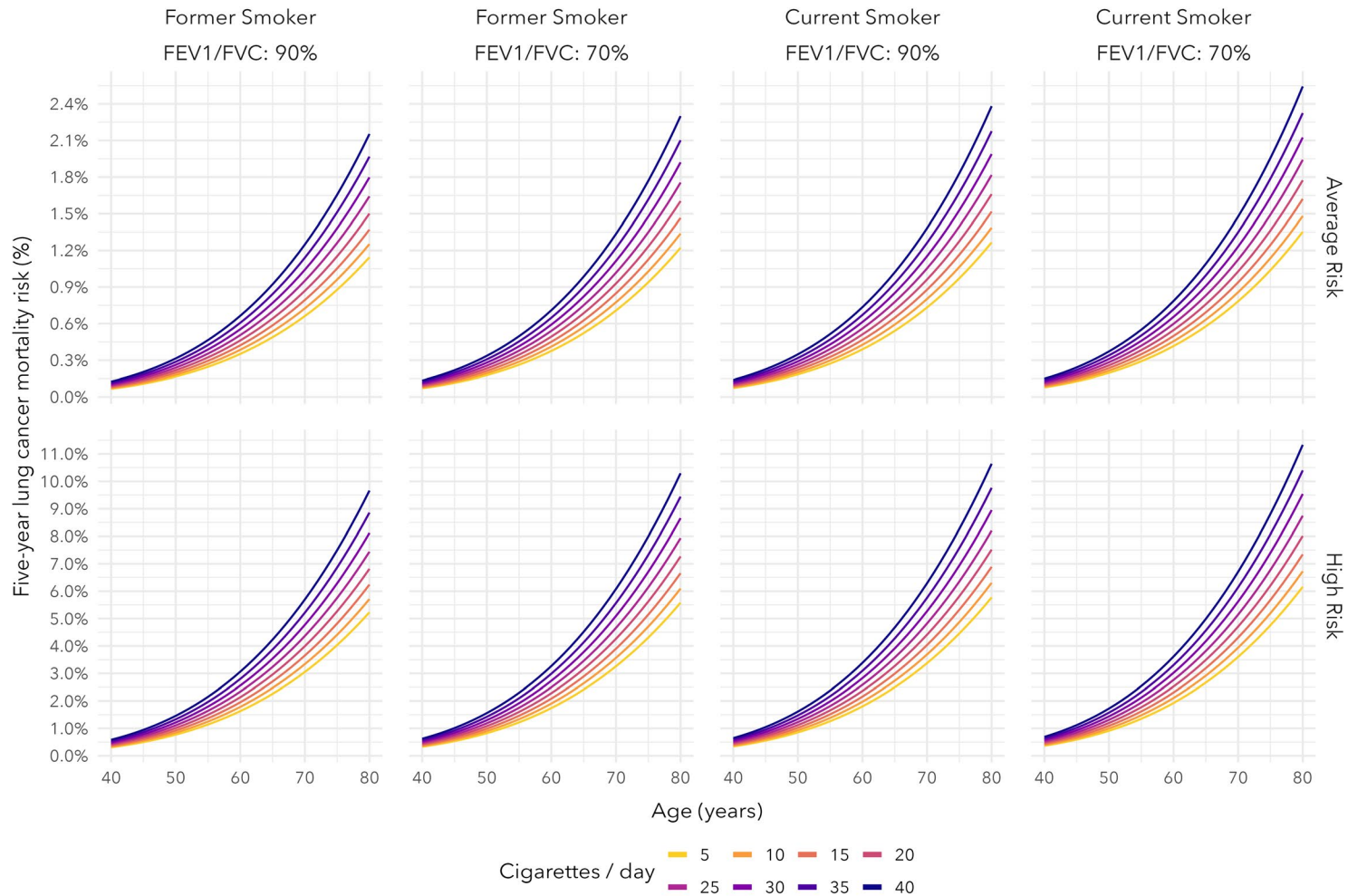


Figure 2: Five-year absolute risk trajectories for lung cancer mortality based on ALARM-ES for current and former smokers for varying smoking intensity (cigarettes per day) and FEV1/FVC, for average and high risk risk profiles. An average risk profile is defined as having the average covariate value for all predictors other than those varied and a high risk profile is defined as having the highest risk covariate value observed in the CKB (based on direction of effect).

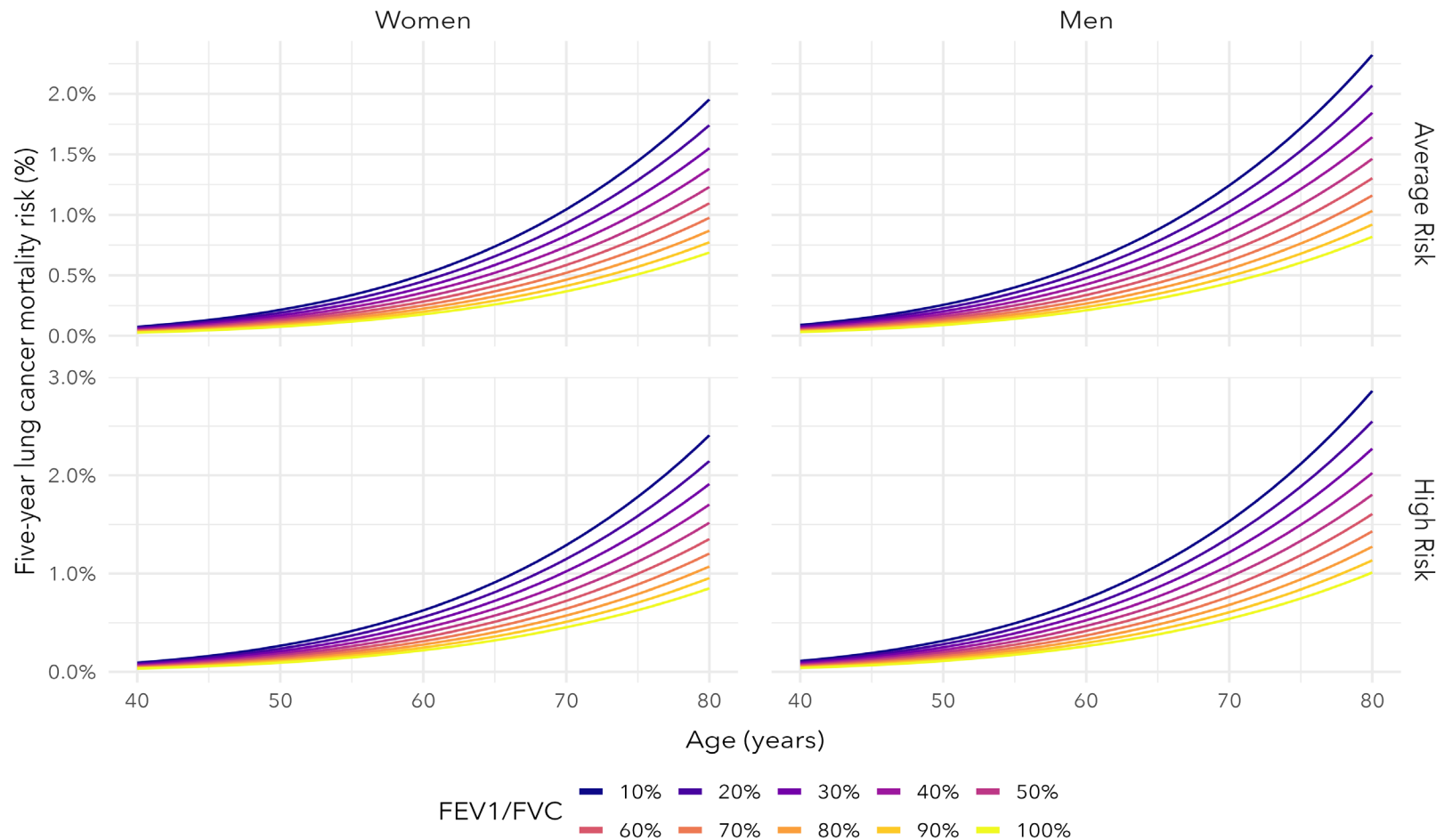


Figure 3: Five-year absolute risk trajectories for lung cancer mortality based on ALARM-NS for never smoker men and women across levels of FEV1/FVC, for average and high risk profiles. An average risk profile is defined as having the average covariate value for all predictors other than those varied and a high risk profile is defined as having the highest risk covariate value observed in the CKB (based on direction of effect).

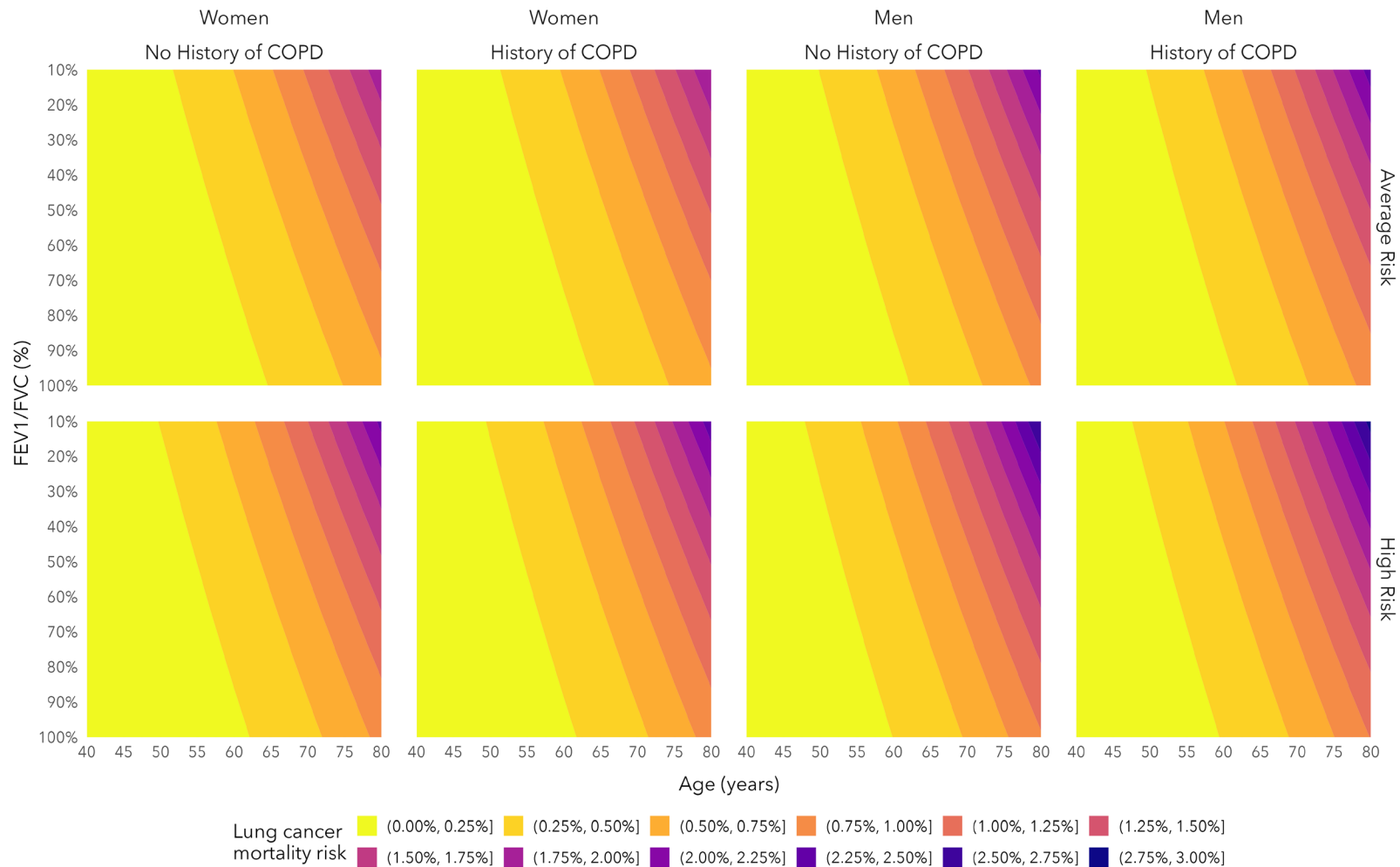


Figure 4: Five-year absolute risk of lung cancer mortality based on ALARM-NS for an average or high risk never smoker for combinations of age and FEV1/FVC. An average risk profile is defined as having the average covariate value for all predictors other than those varied and a high risk profile is defined as having the highest risk covariate value observed in the CKB (based on direction of effect).