
Cost benefit analysis of obesity interventions – technical appendix

HealthLumen on behalf of Nesta



Overview

HealthLumen's simulation consists of two models. The first model uses a regression model which calculates the predictions of risk factor trends over time based on data from rolling cross-sectional studies. The second model performs a monte carlo microsimulation of a virtual population, generated with demographic characteristics matching those of the observed data. The health trajectory of each individual from the population is simulated over time allowing them to contract, survive or die from a set of diseases or injuries related to the analysed risk factors. A detailed description of the two modules is presented below.

Model one: Predictions of risk factors over time

Body mass index (BMI) (as well as other risk factors) is categorised into groups based on the World Health Organization cut-offs [1].

For each RF, let N be the number of categories for a given risk factor, e.g. $N = 4$ for BMI. Let $k = 1, 2, \dots, N$ number these categories and $p_k(t)$ denote the prevalence of individuals with RF values that correspond to the category k at time t . We estimate $p_k(t)$ using multinomial logistic regression model with prevalence of RF category k as the outcome, and time t as a single explanatory variable. For $k < N$, we have

$$\ln\left(\frac{p_k(t)}{p_1(t)}\right) = \beta_0^k + \beta_1^k t$$

The prevalence of the first category is obtained by using the normalisation constraint

$\sum_{k=1}^N p_k(t) = 1$. Solving equation for $p_k(t)$, we obtain

$$p_k(t) = \frac{\exp(\beta_0^k + \beta_1^k t)}{1 + \sum_{k'=1}^N \exp(\beta_0^{k'} + \beta_1^{k'} t)},$$

which respects all constraints on the prevalence values, i.e. normalisation and $[0, 1]$ bounds.

Multinomial logistic regression

Measured data consist of sets of probabilities, with their variances, at specific time values (typically the year of the survey). For any particular time the sum of these probabilities is unity. Typically, such data might be the probabilities of low-risk, medium-risk, or high-risk as they are extracted from the survey data set. Each data point is treated as a normally distributed¹ random variable; together they are a set of

¹ Depending on the circumstances this assumption will be more or less accurate and more or less necessary. In general, it is both extremely useful and accurate. For simple surveys the individual Bayesian prior and posterior probabilities are Beta distributions – the likelihood being binomial. For reasonably large samples, the



N groups (number of years) of K probabilities $\{ \hat{\mu}_k[0, N-1] \}$. For each year the set of K probabilities form a distribution – their sum is equal to unity.

The regression consists of fitting a set of logistic functions $\{ p_k(\mathbf{a}, \mathbf{b}, t) | k \in [0, K-1] \}$ to these data – one function for each k -value. At each time value the sum of these functions is unity. Thus, for example, when measuring obesity in the three states already mentioned, the $k = 0$ regression function represents the probability of being healthy weight over time, $k = 1$ the probability of being overweight and $k = 2$ the probability of being obese.

The regression equations are most easily derived from a familiar least square minimization. In the following equation set the weighted difference between the measured and predicted probabilities is written as S ; the logistic regression functions $p_k(\mathbf{a}, \mathbf{b}; t)$ are chosen to be ratios of sums of exponentials (this is equivalent to modelling the log probability ratios, p_k/p_0 , as linear functions of time).

$$S(\mathbf{a}, \mathbf{b}) = \frac{1}{2} \sum_{k=0}^{K-1} \sum_{i=0}^{N-1} \frac{(p_k(\mathbf{a}, \mathbf{b}; t_i) - \mu_{ki})^2}{\sigma_{ki}^2}$$

$$p_k(\mathbf{a}, \mathbf{b}, t) \equiv \frac{e^{A_k}}{1 + e^{A_0} + \dots + e^{A_{K-1}}}$$

$$\mathbf{a} \equiv (a_0, a_1, \dots, a_{K-1}), \quad \mathbf{b} \equiv (b_0, b_1, \dots, b_{K-1})$$

$$A_0 \equiv 0, \quad A_k \equiv a_k + b_k t$$

The parameters A_0 , a_0 and b_0 are all zero and are used merely to preserve the symmetry of the expressions and their manipulation. For a K -dimensional set of probabilities there will be $2(K-1)$ regression parameters to be determined.

For a given dimension K there are $K-1$ independent functions p_k – the remaining function being determined from the requirement that a complete set of K forms a distribution and sums to unity.

Note that the parameterization ensures that the necessary requirement that each p_k be interpretable as a probability – a real number lying between 0 and 1.

The minimum of the function S is determined from the equations:

approximation of the beta distributions by normal distributions is both legitimate and a practical necessity. For complex, multi-PSU, stratified surveys, it is again assumed that these base probabilities are approximately normally distributed and, again, it is an assumption that makes the analysis tractable.

Depending on the nature of the raw data set it may be possible to use non-parametric statistical methods for this analysis. This is possible for the HSE and GHS data sets of this study but when this has been done the authors can report no discernible difference in the results.



$$\frac{\partial S}{\partial a_j} = \frac{\partial S}{\partial b_j} = 0 \quad \text{for } j=1,2,\dots,k-1$$

noting the relations

$$\frac{\partial p_k}{\partial A_j} = \frac{\partial}{\partial A_j} \left(\frac{e^{A_k}}{1 + e^{A_1} + \dots + e^{A_{k-1}}} \right) = p_k \delta_{kj} - p_k p_j$$

$$\frac{\partial}{\partial a_j} = \frac{\partial}{\partial A_j}$$

$$\frac{\partial}{\partial b_j} = t \frac{\partial}{\partial A_j}$$

The values of the vectors \mathbf{a} , \mathbf{b} that satisfy these equations are denoted $\hat{\mathbf{a}}, \hat{\mathbf{b}}$. They provide the trend lines, $p_k(\hat{\mathbf{a}}, \hat{\mathbf{b}}; t)$, for the separate probabilities. The confidence intervals for the trend lines are derived most easily from the underlying Bayesian analysis of the problem.

Bayesian interpretation

The $2K-2$ regression parameters $\{\mathbf{a}, \mathbf{b}\}$ are regarded as random variables whose posterior distribution is proportional to the function $\exp(-S(\mathbf{a}, \mathbf{b}))$. The maximum likelihood estimate of this probability distribution function, the minimum of the function S , is obtained at the values $\hat{\mathbf{a}}, \hat{\mathbf{b}}$. Other properties of the $(2K-2)$ -dimensional probability distribution function are obtained by first approximating it as a $(2K-2)$ -dimensional normal distribution whose mean is the maximum likelihood estimate. This amounts to expanding the function $S(\mathbf{a}, \mathbf{b})$ in a Taylor series as far as terms quadratic in the differences $(\mathbf{a} - \hat{\mathbf{a}}), (\mathbf{b} - \hat{\mathbf{b}})$ about the maximum likelihood estimate $\hat{\mathbf{S}} \equiv S(\hat{\mathbf{a}}, \hat{\mathbf{b}})$. Hence

$$S(\mathbf{a}, \mathbf{b}) = \frac{1}{2} \sum_{k=0}^{k=K-1} \sum_{i=0}^{i=N-1} \frac{(p_k(\mathbf{a}, \mathbf{b}; t_i) - \mu_{ki})^2}{\sigma_{ki}^2}$$

$$\equiv S(\hat{\mathbf{a}}, \hat{\mathbf{b}}) + \frac{1}{2} (a - \hat{a}, b - \hat{b}) P^{-1} (a - \hat{a}, b - \hat{b}) + \dots$$

$$\approx S(\hat{\mathbf{a}}, \hat{\mathbf{b}}) + \frac{1}{2} \sum_{i,j} (a_i - \hat{a}_i) \frac{\partial^2 \hat{S}}{\partial \hat{a}_i \partial \hat{a}_j} (a_j - \hat{a}_j) + \frac{1}{2} \sum_{i,j} (a_i - \hat{a}_i) \frac{\partial^2 \hat{S}}{\partial \hat{a}_i \partial \hat{b}_j} (b_j - \hat{b}_j) +$$

$$+ \frac{1}{2} \sum_{i,j} (b_i - \hat{b}_i) \frac{\partial^2 \hat{S}}{\partial \hat{b}_i \partial \hat{a}_j} (a_j - \hat{a}_j) + \frac{1}{2} \sum_{i,j} (b_i - \hat{b}_i) \frac{\partial^2 \hat{S}}{\partial \hat{b}_i \partial \hat{b}_j} (b_j - \hat{b}_j)$$

The $(2K-2)$ -dimensional covariance matrix P is the inverse of the appropriate expansion coefficients. This matrix is central to the construction of the confidence limits for the trend lines.



Estimation of the confidence intervals

The logistic regression functions $p_k(t)$ can be approximated as a normally distributed time-varying random variable $N(\hat{p}_k(t), \sigma_k^2(t))$ by expanding p_k about its maximum

likelihood estimate (the trend line) $\hat{p}_k(t) = p(\hat{\mathbf{a}}, \hat{\mathbf{b}}, t)$

$$\begin{aligned} p_k(\mathbf{a}, \mathbf{b}, t) &= p_k(\hat{\mathbf{a}} + \mathbf{a} - \hat{\mathbf{a}}, \hat{\mathbf{b}} + \mathbf{b} - \hat{\mathbf{b}}, t) \\ &= \hat{p}_k(t) + (\nabla_{\hat{\mathbf{a}}}, \nabla_{\hat{\mathbf{b}}}) \hat{p}_k(t) \begin{pmatrix} \mathbf{a} - \hat{\mathbf{a}} \\ \mathbf{b} - \hat{\mathbf{b}} \end{pmatrix} + \dots \end{aligned}$$

Denoting mean values by angled brackets, the variance of p_k is thereby approximated as

$$\begin{aligned} \sigma_k^2(t) &\equiv \left\langle (p_k(\mathbf{a}, \mathbf{b}, t) - \hat{p}_k(t))^2 \right\rangle = (\nabla_{\hat{\mathbf{a}}} \hat{p}_k(t), \nabla_{\hat{\mathbf{b}}} \hat{p}_k(t)) \left\langle \begin{pmatrix} \mathbf{a} - \hat{\mathbf{a}} \\ \mathbf{b} - \hat{\mathbf{b}} \end{pmatrix} \begin{pmatrix} \mathbf{a} - \hat{\mathbf{a}} \\ \mathbf{b} - \hat{\mathbf{b}} \end{pmatrix}^T \right\rangle \times \\ &(\nabla_{\hat{\mathbf{a}}} \hat{p}_k(t), \nabla_{\hat{\mathbf{b}}} \hat{p}_k(t))^T = (\nabla_{\hat{\mathbf{a}}} \hat{p}_k(t), \nabla_{\hat{\mathbf{b}}} \hat{p}_k(t)) P (\nabla_{\hat{\mathbf{a}}} \hat{p}_k(t), \nabla_{\hat{\mathbf{b}}} \hat{p}_k(t))^T \end{aligned}$$

When $K=3$ this equation can be written as the 4-dimensional inner product

$$\sigma_k^2(t) = \begin{pmatrix} \frac{\partial \hat{p}_k(t)}{\partial \hat{a}_1} & \frac{\partial \hat{p}_k(t)}{\partial \hat{a}_2} & \frac{\partial \hat{p}_k(t)}{\partial \hat{b}_1} & \frac{\partial \hat{p}_k(t)}{\partial \hat{b}_2} \end{pmatrix} \begin{bmatrix} P_{aa11} & P_{aa12} & P_{ab11} & P_{ab12} \\ P_{aa21} & P_{aa22} & P_{ab21} & P_{ab22} \\ P_{ba11} & P_{ba12} & P_{bb11} & P_{bb12} \\ P_{ba21} & P_{ba22} & P_{bb21} & P_{bb22} \end{bmatrix} \begin{pmatrix} \frac{\partial \hat{p}_k(t)}{\partial \hat{a}_1} \\ \frac{\partial \hat{p}_k(t)}{\partial \hat{a}_2} \\ \frac{\partial \hat{p}_k(t)}{\partial \hat{b}_1} \\ \frac{\partial \hat{p}_k(t)}{\partial \hat{b}_2} \end{pmatrix}$$

where $P_{cdij} \equiv \left\langle (c_i - \hat{c}_i)(d_j - \hat{d}_j) \right\rangle$. The 95% confidence interval for $p_k(t)$ is centred given as $[\hat{p}_k(t) - 1.96\sigma_k(t), \hat{p}_k(t) + 1.96\sigma_k(t)]$.

Model two: Microsimulation Model

Microsimulation model overview

Simulated people are generated with the correct demographic statistics in the simulation's start-year. In this year women are stochastically allocated the number and years of birth of their children – these are generated from known fertility and mother's age at birth statistics (valid in the start-year). If a woman has children then those children are generated as members of the simulation in the appropriate birth year. The microsimulation is provided with a list of relevant diseases. These diseases used the best available incidence, mortality, survival, relative risk and



prevalence statistics (by age and sex). The virtual population is initialised with diseases by simulating each individual from birth until the start year of the model simulation. It is assumed that a person can die before the model start year. It is assumed that at initialisation the diseases are independent random variables.

During the course of their lives, simulated people can die from one of the diseases caused by the particular risk factor that they might have acquired or from some other cause. The probability that a person of a given age and gender dies from a cause other than the disease are calculated in terms of known death and disease statistics valid in the start-year. It is constant over the course of the simulation. The survival rates from the risk factor-related diseases will change as a consequence of the changing distribution of the risk factor in the population.

The microsimulation incorporates an economic module. The module employs Markov-type simulation of long-term health benefits, health care costs and non-health care related costs of specified interventions. It synthesises and estimates evidence on and cost-utility analysis. The model is used to project the differences in quality-adjusted life years (QALYs), lifetime health-care costs, premature mortality costs and indirect costs as a consequence of interventions over a specified time scale. Outputs can be discounted for any specific discount rate. Incremental cost effectiveness ratios (ICERs) are calculated based on costs per QALY gained

Population module

The population module contains several datasets which can be edited by the end user through a user interface. The population is created in the start-year and propagated forwards in time by allowing females to give birth and also has the ability to incorporate population projections (i.e. migration through minimum arrivals and departures). People within the model can die from specific diseases or from other causes. The <deaths by year by sex by age> file is a necessary input to the model when population projections are being used valid in the start year and usually referred to as the deaths from all causes file. This module is flexible and allows the user to run open and closed populations with no births.

Distributions

Distribution name	symbol	Note
MalesByAgeByYear	$p_m(a)$	Input in year ₀ – probability of a male having age a
FemalesByAgeByYear	$p_f(a)$	Input in year ₀ – probability of a female having age a
BirthsByAgeofMother	$p_b(a)$	Input in year ₀ – conditional probability of a birth at age a the mother gives birth.
NumberOfBirths	$p_l(n)$	l°TFR, Poisson distribution, probability of giving birth to n children
MaleDeathByAge	$p_{Wm}(a)$	Input in year ₀ , probability of a male dying at age a
FemaleDeathByAge	$p_{Wf}(a)$	Input in year ₀ , probability of a female dying at age a



Birth model

Any female in the child-bearing years $\{AgeAtChild.lo, AgeAtChild.hi\}$ is deemed capable of giving birth. The number of children, n , that she has in her life is dictated by the Poisson distribution $p_l(n)$ where the mean of the Poisson distribution is the Total Fertility Rate (TFR) parameter².

The probability that a mother (who does give birth) gives birth to a child at age a is determined from the BirthsByAgeOfMother distribution as $p_b(a)$. For any particular mother the births of multiple children are treated as independent events, so that the probability that a mother who produces N children produces n of them at age a is given as the Binomially distributed variable,

$$p_b(n \text{ at } a | N) = \frac{N!}{n!(N-n)!} (p_b(a))^n (1-p_b(a))^{N-n}$$

The probability that the mother gives birth to n children at age a is

$$p_b(n \text{ at } a) = e^{-\lambda} \sum_{N=n}^{\infty} \frac{\lambda^N}{N!} p_b(n \text{ at } a | N) = e^{-\lambda} \sum_{N=n}^{\infty} \frac{\lambda^N}{n!(N-n)!} (p_b(a))^n (1-p_b(a))^{N-n}$$

Performing the summation in this equation gives the simplifying result that the probability $p_b(n \text{ at } a)$ is itself Poisson distributed with mean parameter $\lambda p_b(a)$,

$$p_b(n \text{ at } a) = e^{-\lambda p_b(a)} \frac{(\lambda p_b(a))^n}{n!} = p_{\lambda p_b(a)}(n)$$

Thus, on average, a mother at age a will produce $\lambda p_b(a)$ children in that year.

The gender of the children³ is determined by the probability $p_{male}=1-p_{female}$. In the baseline model this is taken to be the probability $N_m/(N_m+N_f)$.

Population dynamics

In some year, Y , the population will consist of N_m males and N_f females with their respective age distributions. In the next year, Y' , the numbers will have been depleted by deaths and augmented by the $N_{newborn}$ births. The new, primed, population is determined from the old by the following equation set:

$$N_{newborn} = \lambda N_f \sum_{a=AgeAtChild.lo}^{a=AgeAtChild.hi} p_f(a) (1-p_f(a)) p_b(a)$$

$$N'_m = N_m \sum_{a=1}^{a=Age.hi} p_m(a) (1-p_m(a)) + p_{male} N_{newborn}$$

² This could be made to be time dependent; in the baseline model it is constant.

³ The probability of child gender can be made time dependent.



$$N'_f = N_f \sum_{a=1}^{a=Age.hi} p_f(a)(1-p_f(a)) + p_{female} N_{newborn}$$

$$p'_m(a+1) = \frac{N_m}{N'_m} p_m(a)(1-p_m(a))$$

$$p'_m(a+1) = \frac{N_m}{N'_m} p_m(a)(1-p_{\Omega_m}(a))$$

$$p'_f(a+1) = \frac{N_f}{N'_f} p_f(a)(1-p_{\Omega_f}(a))$$

$$p'_m(0) = \frac{1}{N'_m} p_{male} N_{newborn}$$

$$p'_f(0) = \frac{1}{N'_f} p_{female} N_{newborn}$$

The Population editor' menu item **Population Editor\View\Population dynamics\male** implements these equations and draws projected populations year by year.

Deaths from modelled diseases

The simulation models any number of specified diseases some of which may be fatal. In the start year the simulation's death model uses the diseases' own mortality statistics to adjust the probabilities of death by age and gender. In the start year the net effect is to maintain the same probability of death by age and gender as before; in subsequent years, however, the rates at which people die from modelled diseases will change as modelled risk factors change. This the population dynamics sketched above will be only an approximation to the simulated population's dynamics. The latter will be known only on completion of the simulation.

Multiple population processing

Multiple populations can be used in a simulation provided they are non-overlapping (people cannot belong to both).

In a simulation, Monte Carlo trials are allocated between different populations in proportion to their total person count (malesCount+femalesCount). The idea is to provide a representative sample of the combined population.

In a simulation, a population (pop) is current if the simulated year Y satisfies

$$pop \rightarrow startYear \leq Y \leq pop \rightarrow stopYear$$



Open populations

This model is an *open* population model which allows people to enter and to depart from the population (departure probability $p_d(t)$).

Open population, births and deaths.

In the year y the number of males and females in the population are denoted as $\{N_m(a,y), N_f(a,y)\}$,

And we suppose that they have departure probabilities $\{p_{md}(a,y), p_{fd}(a,y)\}$. The number of new arrivals into each age in the year Y are denoted $\{N_{mArr}(a,y), N_{fArr}(a,y)\}$.

The following analysis applies equally to males and females and we drop the gender suffix. The male and female populations grow according to the recursion relations

$$N(a+1, y+1) = N(a, y)(1 - p_\Omega(a))(1 - p_\delta(a, y)) + N_{Arr}(a, y) \quad (a > 1)$$

$$N(1, y+1) = N_{Newborn}(y)(1 - p_\Omega(0))(1 - p_\delta(0, y)) + N_{Arr}(0, y) \quad (a = 0)$$

The longitudinal modelling of populations having known cross sectional data

Given a set of X-sectional population projections $\{K_m(a,y), K_f(a,y) | 0 \leq a \leq 100; Y_0 \leq y \leq Y_1\}$ (the K- population) the question arises of how to model the lives of individuals within the population (the N-population). In the absence of precise arrival (immigration) and departure (emigration) statistics, many solutions exist. The population is constructed iteratively: given the population in year Y the next year' population is calculated from the known birth and death rates; the departure probabilities and arrival numbers are found by matching with the projected K-population.

Minimum arrival and departure model

The minimum arrival and departure model fixes the modelled N-population in the start year and compensates in subsequent years either by having non-zero departure statistics (if $N > K$) or by importing new people ($K > N$).

From equation :

$$\text{if } N(a, y)(1 - p_\Omega(a)) > K(a+1, y+1)$$

$$(1 - p_\delta(a, y)) = \frac{K(a+1, y+1)}{N(a, y)(1 - p_\Omega(a))} \quad (a > 1)$$

\Rightarrow

$$N(a+1, y+1) = N(a, y)(1 - p_\Omega(a))(1 - p_\delta(a, y)) = K(a+1, y+1) \quad (a > 1)$$

$$\text{if } N(a, y)(1 - p_\Omega(a)) < K(a+1, y+1)$$



$$N_{Arr}(a, y) = K(a + 1, y + 1) - N(a, y)(1 - p_{\Omega}(a)) \quad (a > 1)$$

⇒

$$N(a + 1, y + 1) = N(a, y)(1 - p_{\Omega}(a)) + N_{Arr}(a, y) = K(a + 1, y + 1)$$

The implementation of this model can be arranged using multiple populations – one population for each year of the simulation. The first population consists of the base line model that matches the N and K populations in the start year; subsequent populations contain the corrections (the arrivals, if any in that year). When arrivals enter the simulated population they have a start year corresponding to this population's start year. They usually will have been modelled from birth in the appropriate risk and disease environment. Arrivals are ordinary members of the modelled population – they simply enter the population at times after the simulation-start time. Arrivals carry with them a population identifier.

The numbers of males and females and their ages are known for all populations. Within the micro simulation multiple populations are sampled at a rate proportional to their population size.

Risk factor module

The distribution of risk factors (RF) in the population is estimated using regression analysis stratified by both sex $S = \{\text{male, female}\}$ and age group $A = \text{e.g. } \{0-9, 10-19, \dots, 70-79, 80+\}$. The fitted trends are extrapolated to forecast the distribution of each RF category in the future. For each sex-and-age-group stratum, the set of cross-sectional, time-dependent, discrete distributions $D = \{p_k(t) | k = 1, \dots, N; t > 0\}$, is used to manufacture RF trends for individual members of the population.

Continuous risk factors

In the case of a continuous RF, for each discrete distribution D there is a continuous counterpart. Let β denote the RF value in the continuous scale and let $f(\beta|A, S, t)$ be the probability density function of β for age group A and sex S at time t . Then

$$p_k(t|A, S) = \int_{\beta \in k} f(\beta|A, S, t) d\beta.$$

Equations and both refer to the same quantity. However, equation uses the definition of the probability density function to express the age-and-sex-specific percentage of individuals in RF category k at time t . Equation gives an estimate of this quantity using equation for all $k = 0, \dots, N$. The cumulative distribution function of β is

At time t , a person with sex S belonging to the age group A is said to be on the p -th percentile of this distribution if $F(A, S, t) = p/100$. Given the cross-sectional information from the set of distributions D , it is possible to simulate longitudinal trajectories by forming pseudo-cohorts within the population. A key requirement for



these sets of longitudinal trajectories is that they reproduce the cross-sectional distribution of RF categories for any year with available data. The method adopted here and in the earlier Foresight report [2] is based on the assumption that person's RF value changes throughout their lives in such a way that they always have the same associated percentile rank. As they age, individuals move from one age group to another and their RF value changes so that they have the same percentile rank but of a different RF distribution. In a nutshell, we assume (in accordance with research on the long-term success rate in dieting) that relatively fat people will remain relatively fat and relatively thin people will remain relatively thin. Crucially it meets the important condition that the cross-sectional RF distributions obtained by simulation match the RF distributions of the observed data.

The above procedure can be explained using the example of the alcohol consumption distribution. The alcohol consumption distributions are known for the population stratified by sex and age for all years of the simulation (by extrapolation of fitted model, see equation). A person who is in age group A and who grows ten year older will at some time move into the next age group A' and will have an alcohol consumption that was described first by the distribution $f(\beta|A, S, t)$ and then at the later time t' by the distribution $f(\beta|A', S, t')$. If the alcohol consumption of that individual is on the p -th percentile of the alcohol consumption distribution, their alcohol consumption will change from β to β' so that

$$\beta = F^{-1}\left(\frac{p}{100}|A, S, t\right)$$
$$\beta' = F^{-1}\left(\frac{p}{100}|A', S, t'\right) \Rightarrow \beta' = F^{-1}\left(F(\beta | A, S, t)|A', S, t'\right)$$

Where F^{-1} is the inverse of the cumulative distribution function of β , which we model with a continuous uniform, normal or lognormal distribution (depending on the risk factor) within the RF categories. Equation guarantees that the transformation taking the random variable β to β' ensures the correct cross-sectional distribution at time t' .

The microsimulation first generates individuals from the RF distributions of the set D and, once generated, grows the individual's RF in a way that is also determined by the set D . It is possible to implement equation as a suitably fast algorithm.

Disease module

Disease modelling relies heavily on the sets of incidence, mortality, survival, relative risk and prevalence statistics. The microsimulation uses risk-dependent incidence statistics and these are inferred from the relative risk statistics and the distribution of the risk factor within the population. In the simulation, individuals are assigned a risk factor trajectory giving their personal risk factor history for each year of their lives. Their probability of getting a particular risk factor related disease in a particular year will depend on their risk factor state in that year. The necessary equations are given



below. The microsimulation model has the ability to model discrete multiple stages of a disease.

Once a person has a fatal disease (or diseases) their probability of survival will be controlled by a combination of the disease-survival statistics and the probabilities of dying from other causes. Disease survival statistics are modelled as age and gender dependent exponential distributions.

Relative risks

The reported incidence risks for any disease do not make reference to any underlying risk factor. The microsimulation requires this dependence to be made manifest.

The risk factor dependence of disease incidence has to be inferred from the distribution of the risk factor in the population (here denoted as p); it is a disaggregation process:

Suppose that a is a risk factor state of some risk factor A and denote by $p_A(d|a,a,s)$ the incidence probability for the disease d given the risk state, a , the person's age, a , and gender, s . The relative risk r_A is defined by equation.

$$p_A(d|\alpha,a,s) = \rho_{A|d}(\alpha|a,s) p_A(d|\alpha_0,a,s)$$
$$\rho_{A|d}(\alpha_0|a,s) \equiv 1$$

Where a_0 is the zero risk state (for example, the moderate state for alcohol consumption).

The incidence probabilities, as reported, can be expressed in terms of the equation,

$$p(d|a,s) = \sum_{\alpha} p_A(d|\alpha,a,s) \pi_A(\alpha|a,s)$$
$$= p_A(d|\alpha_0,a,s) \sum_{\alpha} \rho_{A|d}(\alpha|a,s) \pi_A(\alpha|a,s)$$

Combining these equations allows the conditional incidence probabilities to be written in terms of known quantities

$$p(d|\alpha,a,s) = \rho_{A|d}(\alpha|a,s) \frac{p(d|a,s)}{\sum_{\beta} \rho_{A|d}(\beta|a,s) \pi_A(\beta|a,s)}$$

Previous to any series of Monte Carlo trials the microsimulation program pre-processes the set of diseases and stores the *calibrated* incidence statistics $p_A(d|a_0, a, s)$.

For each scenario the incidence statistics are calibrated against the baseline trends.



Model output module

Cross-sectional outputs (epidemiological and economic) per 100,000 of the population are computed for each year of the simulation.

Epidemiological and economic outputs

A range of different epidemiological outputs are produced by the model including:

- Cumulative incidence rates
- Quality Adjusted Life Years (QALY)
- Incremental cost-effectiveness ratio (ICER)
- Net monetary benefit (NMB)

The QALY outputs can be discounted if required and this can be defined by the user at the start of a modelling project. The discounting rate each year ($Discount(year)$) was calculated as shown in equation (1.35).

$$Discount(year) = \frac{1}{(1 + R)^{year - year_{start}}} \quad (0.33)$$

Where, year start refers to the start year of the modelling which is 2019 in this study and R is the annual discount rate. The following sections describe some of these outputs in more detail.

ICER:

$$ICER = \frac{\text{incremental costs}}{QALYs \text{ gained}} \quad (0.34)$$

NMB:

$$NMB = QALYs \text{ gained} \times \text{£}20,000 - \text{incremental costs} \quad (0.35)$$

Where *incremental costs* is the additional costs incurred as the result of the policy, and £20,000 is the willingness-to-pay / cost-effectiveness threshold for the UK.

Confidence intervals

95% confidence intervals for microsimulation outputs were calculated at the single age / sex level using the formula:

$$CI_{x|a,s} = \pm 1.96 \cdot SE_{x|a,s} \cdot W_{a,s} \quad (0.36)$$

Where:

- x is an outcome measure or event, e.g. prevalence of a type 2 diabetes;



- a is the group's age, in years, from 0 to 110;
- s is the group's sex, of *male* and *female*;
- $CI_{x|a,s}$ is the confidence interval for x given age a and sex s ;
- $W_{a,s}$ is the population weight for age a and sex s (i.e., the number of people in the population with that age and sex for that year); and
- $SE_{x|a,s}$ is the standard error of the estimate for event x , age a and sex s , as estimated by:

$$SE_{x|a,s} = \sqrt{\frac{p_{x|a,s}(1-p_{x|a,s})}{N_{a,s}}} \quad (0.37)$$

Where $N_{a,s}$ is the number of microsimulation trials for age a and sex s , and $p_{x|a,s}$ is the proportion of this group in which the event/outcome has been recorded.

Error propagation

When adding or subtracting confidence intervals – for example, when summing all ages and sexes to achieve population-level estimates, or when calculating incremental costs – the total confidence interval was calculated as follows:

$$\sqrt{\sum_{i=1}^n CI_i^2} \quad (0.38)$$

Where CI_i is one of n confidence intervals..

For calculation of the ICER CI we used the following equation:

$$\sqrt{\left(\frac{CI_c}{x_c}\right)^2 + \left(\frac{CI_q}{x_q}\right)^2} \cdot ICER \quad (0.39)$$

Where:

- CI_c and CI_q are the confidence intervals for incremental costs and QALYs gained, respectively;
- x_c and x_q are the estimates for incremental costs and QALYs gained, respectively; and
- $ICER$ is the incremental cost-effectiveness ratio, calculated as described above.

Modelling policy scenarios

Specific assumptions were required for modelling each obesity policy. Nesta provided modelled data quantifying the impact of each policy on BMI and this



methodology is described elsewhere. This data was ingested into the HealthLumen microsimulation model.

The majority of policies followed the same intervention structure. During the first year of the simulation, an individual's eligibility for the intervention was assessed based on their age and BMI. For policies that affect children, to account for different obesity classifications based on age and sex, children were eligible if their BMI was \geq the 85th percentile according to UK90 reference data. If an individual was eligible, then the individual received the effect of the intervention in the following year, and any subsequent years thereafter, and experienced a reduction in their BMI. The BMI percentile is then recalculated using the new BMI value of that individual. For policies impacting a subset of children only, the impacts are applied only to children of that age, not when they move out of that age category.

For policies impacting both children and adults, the effect of the intervention for adults is applied once the child turns 18.

For the policy "Everyone with a BMI of 30 or above is offered a total diet replacement (TDR)", which only affects adults, assessing eligibility remains the same. In this policy, 40% of individuals eligible were offered TDR, of which 13% agreed to undergo TDR. This was then translated into the model as 5.2% receiving TDR. Those that were offered TDR were not offered it in any subsequent years, and as such did not have another chance to be intervened upon. Individuals regain weight following the intervention. Individuals who underwent the intervention had a 33.33% chance of experiencing type 2 diabetes remission in the first year after the intervention [3]. Any individual who achieves type 2 diabetes remission can contract type 2 diabetes again based on disease incidence statistics after the first-year post-intervention. Eligible individuals who do not receive the intervention in that year are eligible for the intervention in the following year, given their BMI remains at 30 kg/m² or above.

For the policy "Introducing universal free school meals for all primary school children during term time", in 2020, all children in school aged 5 to 12 were eligible to receive the intervention each year. Children with a BMI higher than the 95th percentile using the UK90 thresholds had a probability of their BMI being impacted on and reduced to the 94th percentile in that year, after which they follow trends for children on that percentile over time. The probability of being impacted by the intervention varied by age: eligible children aged 5 to 7 had an 8.7% chance, eligible children aged 8 to 10 had a 6.5% chance, and eligible children aged 11 and 12 had a 2.8% chance of having their BMI percentile capped. Children with a BMI above 95th percentile who turned 5 in subsequent years (2021 to 2024) were also eligible to have an 8.7% probability of being impacted by the intervention in that year.

For the policy "Invest a further £500 million over 5 years in local authorities to plan and deliver active transport through Active Travel England (or equivalent in DAs)", a policy that affects adults. In this policy, eligible individuals had a 1.65% chance of receiving the effect of the intervention in the following year, and any subsequent years thereafter, and experienced a reduction in their BMI.



For the policy “Require 60 minutes of daily physical activity for school children” from 2020 to 2024, all children in school aged 5 to 17 with a BMI over the 85th percentile using the UK90 thresholds were eligible to receive the intervention, as long as there were available spaces in the programme. Once all spaces in the programme quota were met, eligible children would no longer be impacted on by the intervention. In the start year of 2019 the quota was set at 1,744,786 places, and the number of available places each year was determined by the number of eligible children (depending on age and BMI). Children with a BMI above 95th percentile who turned 5 in each modelling year were also eligible to be impacted by the intervention in the next years, depending on availability of programme spaces.

For the policy “Extend access to pharmacological interventions by providing an extra £500 million per year of ring-fenced funding to provide increased access to NICE recommended weight loss treatments (Liraglutide and Semaglutide)”, a policy that only affects adults, eligibility was based on age, BMI, and ethnicity. In the scenario, adult individuals had a 13% chance of being Black or Asian (Census, 2021). Individuals were eligible for the intervention if they were ≥ 18 years old, and met one of the following criteria: 1) if they were non-White, and had a BMI of over 27.5 kg/m^2 ; or 2) if their BMI was 30 kg/m^2 or above. Of those eligible, approx. 148,810 received the intervention each year. BMI reduction was applied in the same year. Individuals with type 2 diabetes had a 41.2% remission rate in the first year after the intervention. This was based on a phase 4 observational study, where individuals with HbA1c levels of $< 7.0\%$ at 33-44 weeks were assumed to be in type 2 diabetes remission, where we have assumed the outcome of the study is an annual remission rate [4]. As clinical remission is $< 6.5\%$ HbA1c levels, we may therefore be overestimating remission. Any individual who achieves type 2 diabetes remission can contract type 2 diabetes again based on disease incidence statistics after the first-year post-intervention.

For policy “Extend access to pharmacotherapy so that approximately 3 million more people ($\text{BMI} \geq 30$) receive Semaglutide each year”. After the first year, each year individuals were eligible for the intervention if they were ≥ 18 years old and their BMI was 30 kg/m^2 or above. Of those eligible the BMI reduction was applied in the same year, but only up to 3,000,000 could receive the intervention each year. In the final year of the simulation (2024), only 1,030,394 individuals could receive the intervention, to align with Nesta’s modelling. Individuals with type 2 diabetes had a 41.2% remission rate in the first year after the intervention, as above for policy 24. Any individual who achieved type 2 diabetes remission could contract type 2 diabetes again based on disease incidence statistics after the first-year post-intervention.

For the policy “Expand NHS provision of bariatric surgery to individuals with BMI ≥ 35 with a pre-existing condition (specifically double the amount of people receiving surgery from approximately 6,500 per year to 13,000 per year)”, a policy that only affects adults, eligibility was based on whether the individual had type 2 diabetes or cardiovascular disease. If the individual had one of these two diseases, they were eligible to receive the intervention. Any individual who received the intervention had a 54% chance each year between 2020 and 2024 to achieve remission from hypertension within the intervention, and then a 58.8%, 57.7%, 56.7%, 55.8% an



54.7% for each consecutive year within the intervention to achieve remission from type 2 diabetes [5-7]. If an individual achieved remission for either disease, then they were no longer considered a prevalent disease case. The number of bariatric surgeries within the model that individuals could receive was 6,500 per year, to model the additional 6,500 surgeries each year required to double the amount of people receiving bariatric surgery. If the quota of surgeries was reached, no other individuals would receive the intervention, and the microsimulation would progress to the next year.

For the policy “Provide £85 million of funding per year for increased roll-out of family-based programmes to the local authorities with the highest childhood obesity rates”, a policy that only affects children, the eligibility assessment is the same as the general policy method (aged between 5 and 18, and a BMI percentile of the 85th or higher), with two additional criterium: the child must live in QIMD 4 or 5, and have a parent living with excess weight. For children that were eligible, only 462,187 a year could receive the intervention (number provided by Nesta: £85 million / £320 per family * 1.74 children per family), and the effect of the intervention was applied in the same year of the simulation. In this scenario, individuals could only receive the intervention once.

For the policy “Allocate £100 million per year to fund a programme of financial incentives to improve health behaviours in local authorities with the highest obesity rates”, a policy that affects adults only, the eligibility assessment is the same as the general policy method (age 18 to 100 and BMI \geq 25). However, for adults that were eligible, only 465,116 could receive the effect of intervention.



References

1. WHO. *A healthy lifestyle - WHO recommendations*. 2010. 6th May 2010 [cited 2024 18th July]; Available from: <https://www.who.int/europe/news-room/fact-sheets/item/a-healthy-lifestyle---who-recommendations>.
2. Butland, B., et al., *Tackling obesities: future choices-project report*. Vol. 10. 2007: Citeseer.
3. Lean, M.E.J., et al., *Primary care-led weight management for remission of type 2 diabetes (DiRECT): an open-label, cluster-randomised trial*. *The Lancet*, 2018. **391**(10120): p. 541-551.
4. van Houtum, W., et al., *Real-World Use of Oral Semaglutide in Adults with Type 2 Diabetes in the PIONEER REAL Netherlands Multicentre, Prospective, Observational Study*. *Diabetes Therapy*, 2024.
5. Canakis, A., et al., *Type 2 Diabetes Remission After Bariatric Surgery and Its Impact on Healthcare Costs*. *Obesity Surgery*, 2023. **33**(12): p. 3806-3813.
6. Moriconi, D., M. Nannipieri, and E. Rebelos, *Bariatric surgery to treat hypertension*. *Hypertension Research*, 2023. **46**(5): p. 1341-1343.
7. Toolabi, K., et al., *Comparison of Laparoscopic Roux-en-Y Gastric Bypass and Laparoscopic Sleeve Gastrectomy on Weight Loss, Weight Regain, and Remission of Comorbidities: A 5 Years of Follow-up Study*. *Obesity Surgery*, 2020. **30**(2): p. 440-445.