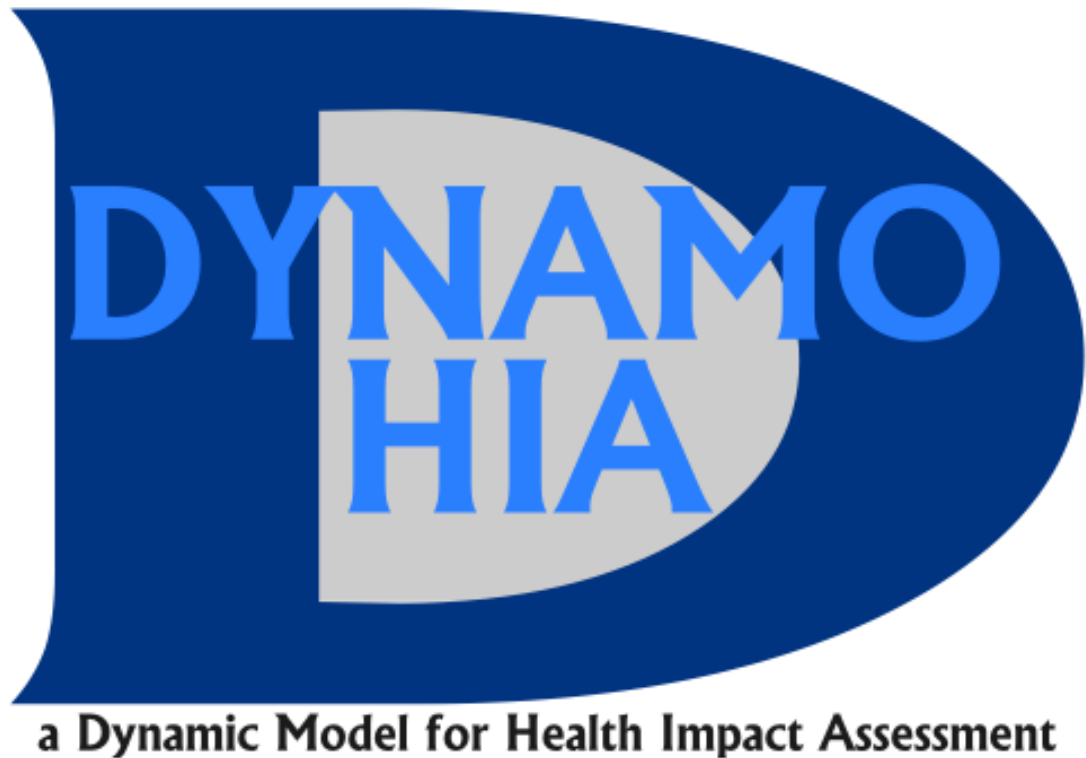


Detailed description of DYNAMO-HIA calculations

Version 1.0



H.C. Boshuizen.

October 2010

This investigation has been performed partially by order and for the account of DG-SANCO, within the framework of grant agreement 2006116

Table of Contents

1	Introduction and general description of the DYNAMO-HIA Model4	
2	Description of the model states	6
3	Estimation of transition rates and occupancy rates of the initial population 9	
3.1	General overview	9
3.1.1	Converting transition rates into transition probabilities	9
3.1.2	Overview of transition rates in the model	11
3.2	Estimation of disease probabilities in the initial population	13
3.2.1	Splitting the cancer prevalence into “cured” and “non-cured” prevalence	14
3.2.2	Calculating excess mortality rate from median survival	15
3.2.3	Calculating the probability of disease in an individual	16
3.3	Estimation of the parameters of the transition rates.	22
3.3.1	Estimation of baseline incidence rates	24
3.3.2	Estimation of baseline fatal disease rates	26
3.3.3	Estimation of attributable mortality rates and other cause mortality	27
3.3.4	Estimation of baseline other mortality and the relative risks for other mortality	31
3.3.5	Estimation of Non user-specified risk factor transition rates (nett transition rates)	34
4	Generation of Initial Population	39
4.1	Population characteristics	39
4.2	Generation of characteristics for the reference population	40
4.3	Generation of characteristics for scenario populations	43
5	Description of simulation module	46
5.1.1	Synchronising of scenarios	46
5.1.2	Adding newborns to the population	47
5.1.3	Update rules for characteristics of the simulated persons.	47
6	Post-processing module and output generated	56
6.1	Weighting procedure	56
6.2	Content of DYNAMO-output files	57
6.3	Calculation of integral measures of health.	59
	Reference List	62

1 Introduction and general description of the DYNAMO-HIA Model

This document describes in detail the calculations that are carried out in the DYNAMO-HIA model. This document assumes that the reader is already familiar with the general aim and design of the DYNAMO-HIA -model. For a general introduction, we refer to Lhachimi et al.[1]. and the DYNAMO-HIA user manual⁷, and for a more concise description of the calculation methods to Boshuizen et al.[2].

The DYNAMO-HIA model is in the form of a “partial” micro-simulation model. The modelling takes place in the following steps:

1. An initial population of individuals is generated, based on the risk factor exposure data that are given to the model, like percentage of smokers or the distribution of BMI. Also the probability of disease is calculated for each simulated individual in this initial population, based on the data on disease prevalence rates that are given to the model.
2. The risk factor histories of these individuals are simulated over time (simulation module), both in the current situation (business-as-usual scenario, which we will further call the **reference scenario**), and under one or more policy scenarios (further referred to as **alternative scenarios**).
3. From these risk factor histories, and the disease probabilities in the initial state, the probabilities on diseases and mortality for each simulated individual during simulated time are calculated, both in the current situation, and under the policy scenario's.

We call this partial micro-simulation because micro-simulation is used to generate risk factor histories, but not for generating disease histories. The probabilities of disease of each individual are calculated using deterministic methods. This choice is made because incidences on diseases are often small, so large numbers of simulated persons would be needed for a sufficiently accurate simulation. Within the Dynamo model, therefore, using deterministic calculations of disease probabilities is more efficient.

Both the simulation of risk factor histories and the calculation of disease probabilities is based on a Markov Model. A Markov Model is described by states, and transition probabilities between states. Chapter 2 will describe the states used in the model, while chapter 3 will describe the transition rates that are used and the way the parameters of these transition rates are estimated from the input data. Chapter 4

will describe the generation of the initial population for the simulation, chapter 5 the simulation, and chapter 6 the post-processing calculations, in which the results from the simulated population are used to calculate 1) the number of persons in the real population with and without diseases, and 2) integrated measures as life expectancies, DALEs and life expectancies without disease.

2 Description of the model states

The state of an individual in the DYNAMO-HIA model is determined by the (joint) value of the following characteristics:

1. **Age**
2. **Gender**
3. **Risk factor status**

The risk factor in the model can be one of three forms:

- A risk factor with a continuous distribution (for instance BMI). In this case the model needs the following input: the shape of the distribution (normal or lognormal), mean of the distribution, the standard deviation of the distribution and in case of the lognormal distribution also the skewness of the distribution. We will further call this **continuous risk factor**.
- A risk factor in classes (for instance: smokers / former smokers / never smokers). In this case the model needs as input: the percentage in each class. A maximum of 10 classes can be entered in the program. We will further call this **categorical risk factor**.
- A risk factor in classes (as above), but with the addition that for one class also the duration of being in this class is of interest (for instance in case of smoking the duration of having stopped smoking is of interest). We will further call this **compound risk factor**. At the start of simulation the duration can have the values 0 to 19.

Only one of these risk factor types can be entered in the model. However, the user could use the 10 classes available to define “combination risk factors” like “smoking alcohol drinkers”, “non smoking alcohol drinkers” “smoking abstainers” and “non-smoking abstainers” .

4. **Disease probabilities**

The disease states in the underlying conceptual model for each disease are either 0 (disease not present) or 1 (disease present). In the DYNAMO-HIA model, however, the “state”, that is, the value stored during simulation for this variable is not the absence or presence of disease, but the probability of having a disease. This approach, using probabilities, is preferred over the method in which an explicit disease state is simulated for each person, as many diseases are rare, so extremely large numbers of persons need to be simulated to obtain sufficient precision.

We will use the following terminology with regard to diseases:

Independent diseases

With independent we mean that the incidence of these diseases does not depend on the presence or absence of other diseases, given the risk factor history (conditional independence).

Dependent diseases

These diseases can depend on the presence/absence of an intermediate disease (but in the DYNAMO-HIA model they may not depend on the presence/absence of another dependent disease)

Intermediate diseases or causal diseases

These are independent diseases that are themselves a risk factor for a dependent disease: such a disease has an intermediate role in the causal pathway between risk factors and the dependent disease. In the DYNAMO-HIA model a dependent disease can not be an intermediate disease, but it can have several dependent diseases. In other words, the causal network of diseases can only have two layers.

The diseases will be assigned to **disease clusters**. A dependent disease is always in the same cluster as the intermediate diseases that are its risk factors. Independent diseases that are not intermediate diseases form a cluster of one (also referred to as: **single disease**).

For each disease cluster, a person’s state will be characterized by the set of probabilities of having each possible combination of diseases within the cluster (conditional on being still alive). For instance, in case of a cluster with diseases A, B, and C, the person will be characterized by the probability of having no disease, of having disease A only, of having disease B only, of having disease C only, of having both disease A and B, of having both disease A and C, of having both disease B and C and of having all three diseases (A, B, and C). In the case of an independent disease that is not an intermediate disease (a cluster

of one), this will simply be the probability of having this disease (conditional on being alive), or not having this disease. As is shown in Boshuizen et al.[2], the overall probability of all combinations of all diseases can be calculated by multiplying the probabilities from the different clusters.

This can also be seen as constructing a multi-state life table for a group persons that have the same risk factor history as the simulated person.

Furthermore, for diseases for which the user specifies a cured fraction (an option added for better simulation of cancers) the probability of disease is separated in a “cured” fraction and a fraction that in time will die of the cancer. The cured fraction is assumed to be independent of risk factor status or the presence of other diseases.

Another option is to define a “fatal Incidence” to model diseases that have an elevated mortality right after contracting the disease (e.g. stroke) but after some time a lower constant excess mortality.

5. Survival probability

Like for diseases, mortality can be a rare event for some ages and thus simulation is inefficient. Therefore for each person at each moment in time, the state that is kept track of in the simulation is the survival probability.

3 Estimation of transition rates and occupancy rates of the initial population

3.1 General overview

During simulation, the value of the four types of characteristics (age, risk factor, disease probability, survival probability) is updated at one year intervals. The characteristic sex is constant, so it does not need updating. The basis for the update are the transition rates between the states of the underlying Markov model. With the underlying Markov model we mean the Markov model that has age, gender, risk factor state, disease state and survival state as its states, rather than age, gender, risk factor state, disease probability and survival probability. The underlying Markov Model has the following transitions:

- transitions of age (simply implying that age increases with one year at each update)
- transition rates between different risk factor states
- transition rates from not having disease A to having disease A
- transition rates from being alive to being dead

Note that no transition rate is included from having the disease to not having the disease, so the model does not allow for remission.

3.1.1 Converting transition rates into transition probabilities

The time-step used in the Dynamo model is 1 year, although in principle it could have been any other period. Most input to the DYNAMO-HIA model is in the form of transition rates (for instance, disease incidence and mortality), where a rate is the change in an infinitely small time period. When updating the characteristics of each simulated person, all transition rates have to be converted into 1-year transition probabilities. Exceptions are the transitions between risk factor states, where the model expects one-year transition probabilities as input rather than transition rates, and age, where no user input is required.

If only transitions out of a state are possible, conversion from rates to 1 year probabilities can be done by using the formula:

$$P(1 \text{ year}) = 1 - e^{-\text{rate} * 1 \text{ year}} \quad (1)$$

where $P(1 \text{ year})$ is the 1 year transition probability and rate is the transition rate. However, in cases where there are also transitions into the state, things are more complicated.

A full analytical solution can be calculated by:

$$\mathbf{P} = e^{\mathbf{A} * \text{timestep}} \quad (2)$$

where \mathbf{P} is the matrix of transition probabilities for the states involved and \mathbf{A} is a matrix of transition rates.

Calculating a matrix exponential is an old numerical challenge[3]. Gallivan and others suggested an algorithm specifically for this type of models[4]. DYNAMO-HIA uses the latter algorithm to solve the equation above.

Furthermore, the computer time required for calculating a matrix exponential increases exponentially with the size of the matrixes involved. In DYNAMO-HIA, we do not include risk factors as a state in this type of calculation, but repeat the calculations for each risk factor state separately, as this makes definition of risk factors and especially of risk factor histories much more flexible. However, the size of the matrix still doubles with every disease added to the model. In order to keep both the dimension of the matrix small, and the calculations transparent, DYNAMO-HIA does not carry out calculations for the entire disease-space, but divides the transition matrix in sections that can be updated independently: Groups of diseases that are independent of other diseases (conditional on risk factor status) (**disease clusters**) are updated as a group, and calculation of the transition matrix is performed within this group. The updates of the disease states per cluster, and the survival due to mortality from other causes of death (that is, from diseases not included in the model) can be used to calculate every disease state, as well as the overall survival at each moment in time. A more detailed description and the justification of this approach can be found in the technical paper describing DYNAMO-HIA[2].

3.1.2 Overview of transition rates in the model

Table 1 describes the transition rates in the DYNAMO-HIA model. The equations given for the transition rates in this table contain constants such as “baseline incidence”. We will further refer to these constants as **(model)parameters**. Some of these constants have to be given by the users (like relative risks on diseases) and are used directly as given, while others are calculated in a parameter estimation module from other data given by the user. In section 3.3 we will describe the way in which we estimate these parameters from the input data.

The element “RR(riskfactor)” in table 1 stands for a function that differs for different types of the risk factors. For the categorical risk factors RR(riskfactor) is equal to the values given in the input. For the other types of risk factors they are calculated as follows:

Continuous risk factors:

$$RR(\text{risk-value}) = RR(\text{per unit})^{(\text{risk-value} - \text{reference value})} \quad (3)$$

Here the reference value is an arbitrary level of the risk factor for which the RR is set to 1. The user can choose this value. However, choosing a different reference value does not influence results.

Compound risk factors:

For the classes without a duration value: RR(riskfactor) is given directly by the user (in the same way as for the categorical risk factor) and for the particular class that models duration:

$$RR = RR_{\text{end}} + (RR_{\text{begin}} - RR_{\text{end}}) \exp(-\alpha * \text{duration in years}) \quad (4)$$

Here

RR_{end} (relative risk at the end of the duration)

RR_{begin} (relative risk at the beginning of the duration), and

α (steepness parameter for the decline) have to be given directly by the user.

Table 1: transition rates in the DYNAMO-HIA model of person i. All transition rate formulas are on the level of the individual, so all probabilities refer to the probability for this individual, not to population probabilities

Type of transition	Quantification of this transition rate
Risk factor state → other risk factor state	<p>Categorical: Given by user¹⁾ or calculated from input prevalence of the risk factor assuming that risk factor prevalence over age groups does not change in time (see 3.3.5).</p> <p>For compound variables, the duration in the duration class is increase by 1 year at each time step, for the none duration classes like above</p> <p>Continuous (normally distributed²⁾): $\text{new value}_i = \text{old value}_i + \text{drift} + \text{stdDrift} * \varepsilon$, with Drift =: either user specified or calculated from input prevalence of the risk factor assuming that risk factor prevalence over age groups does not change in time (see 3.3.5). stdDrift = calculated from input prevalence of the risk factor assuming that the increase in variance of the risk factor prevalence over age groups does not change in time (see 3.3.5). ε = randomly drawn value from a standard Normal distribution</p>
Without disease → with disease (independent diseases)	<p>$\text{Baseline incidence}_d * RR_d(\text{riskfactor}_i)$.</p> <p>$RR_d(\text{riskfactor}_i)$ = relative risk on disease d due to risk factor status</p> <p>$\text{Baseline incidence}_d$ = incidence rate of disease d in those with $RR_d(\text{riskfactor}_i) = 1$</p> <p>$d$ = index of disease</p>
Without disease → with disease (dependent diseases)	<p>$\text{Baseline incidence}_d * RR_d(\text{riskfactor}_i) * \sum_{combi} P_i(combi) RR_d(combi)$.</p> <p>$RR_d(combi)$ = relative risk on disease d of a particular combination of intermediate diseases</p> <p>$P_i(combi)$ = probability of the combination of intermediate disease</p> <p>$combi$ = index of each possible combination²⁾ of intermediate diseases</p>
→ death	<p>$RR_{oc}(\text{riskfactor}_i) * \text{Baseline other mortality} + \sum_d AM(d) P_{nc,i}(d)$</p> <p>$+ \sum_d \{RR_d(\text{riskfactor}_i) \sum_{combi} RR_d(combi) \text{Baseline} F(d)\}$</p> <p>$RR_{oc}(\text{riskfactor}_i)$ = relative risk on other mortality due to risk factor status</p> <p>Baseline other mortality = Other mortality rate in those with $RR_{oc}(\text{riskfactor}_i) = 1$</p> <p>$AM(d)$ = Mortality attributable to disease d (Attributable Mortality)</p> <p>$P_{nc,i}(d)$ = probability of having disease d (excluding cured fraction)</p> <p>Baseline F(d) = Acutely fatal incidence rate of disease d (Fatal) for those with $RR(\text{riskfactor}_i) = 1$</p>

1) Here the input asked or calculated is not transition **rates**, but one-year **transition probabilities**. This means that the conversion from rates to one-year probabilities is not necessary

2) For lognormal distributed risk factors, see see 3.3.5

3) including the combination “having none of the diseases”

The mortality as defined in table 1 can accommodate three disease-mortality processes:

1. For a chronic disease, for which the excess mortality (defined as the difference in mortality rate between those with the disease and those without the disease) only depends on age and gender, but not on how long one has the disease ($\text{baseline}F(d)=0$).
2. For a partly acutely fatal disease: This are diseases (like myocardial infarction or stroke) that occur as a distinct event, that has a very high mortality rate immediately after the event, while those who survive this critical period have a lower mortality, although still higher than the general population. The excess mortality rate after the critical period, like the excess mortality in the first disease process, depends only on age and gender, and not on the duration of the disease ($\text{baseline}F(d)>0$).
3. For a disease with a cured fraction in which the excess mortality is zero in a part of the patients (cured fraction), while the excess mortality does not depend on duration of disease in the others (Which is equivalent to assuming an exponential survival model with a cured fraction). As the part of the patients that are “cured” can only be identified in retrospect (after all patients that have not been cured have died), it is not realistic to model “cured” with remission as with chronic diseases it is not known that they are cured but only a mortality that does not differ from non-diseases is observed. Hence, these diseases are modelled by regarding them as two separate diseases, and putting $A_m(d)$ to zero for the cured disease. Combinations of fatal fractions and cured fractions are not allowed in the DYNAMO model.

3.2 Estimation of disease probabilities in the initial population

Before giving the methods used to estimate the parameters of the transition rates between states, we first describe the methods used to estimate the disease prevalence rates in the initial population, as these are used to estimate several of these parameters.

The input of the model consists of disease prevalence on the population level. These should be translated into prevalence rates at the individual level. For instance, the model input consists of the prevalence of diabetes, and the prevalence of coronary heart disease on the population level. To generate an initial population, one has to know also the prevalence of having diabetes and coronary heart disease simultaneous, and also how the prevalence rates of diseases are related to the risk factors in the model. This needs to be known before starting the estimation of the model parameters (the parameters of the

transition rates), because several estimation procedures assume that the prevalence rate of diseases on the individual level is known.

The individual prevalence rates can not be inferred from the input data as given by the user to the DYNAMO-HIA model without further assumptions. Here we will throughout make the assumption that the prevalence odds ratio of a disease in two groups is equal to the ratio of the incidence rates of this disease in those groups[5, 6]. This is only true in the case that incidence relative risks and excess mortality are constant (do not depend on age within the period that a disease lasts), which clearly will not always be the case. However, we believe this assumption to be superior to the assumption of independence, which is implicitly used in many other models.

Secondly, for cancers the prevalence needs to be split in the prevalence of “cured” patients, and non-cured patients.

3.2.1 Splitting the cancer prevalence into “cured” and “non-cured” prevalence

For cancers the prevalence needs to be split in the prevalence of “cured” patients, and non-cured patients. In a period with constant incidence and a disease without excess mortality (as cured cancer), the relation between the prevalence at age t_2 and at age t_1 is:

$$p(t_2) = p(t_1)(1 - e^{-incidence*(t_2-t_1)}) \quad (5)$$

One option is to use this formula to calculate the prevalence of cured patients at age $t+1$ from the prevalence at age t , using as incidence at that age the incidence times the cured fraction, starting at age 0 with a prevalence of 0.

This method assumes that there are no cohort effects. In reality, however, it might be the case, due to medical advances, that the current “cured fraction” is larger as it was in the past. In that case this method implies that a percentage of cured cases is used for the past that is too large. This then means that a too large part of the current prevalent cases is considered cured, and as cured patients will remain alive (apart from being subject to general mortality) during simulation, the projected prevalence rates will be too high. The second option is to calculate the prevalence of non-cured patients at age $t+1$ ($p(t+1)$) from the prevalence of non-cured patients at age t ($p(t)$) and the incidence of non-cured cancer $inc(t)$ ($= incidence * (1 - cured\ fraction)$) using:

$$p(t+1) = \frac{(p(t)M(t) - Inc(t))e^{(Inc(t)-M(t))} + Inc(t)(1-p(t))}{(p(t)M(t) - Inc(t))e^{(Inc(t)-M(t))} + M(t)(1-p(t))} \quad (6)$$

where $M(t)$ is the excess mortality rate at age t in those with (non-cured) cancer and $Inc(t)$ the incidence at age t . $M(t)$ is given by the user, or it is calculated from the median survival.

This approach is less sensitive to cohort effects as the first option, as only cohort effects during the “dying period” after incidence are relevant. For most cancers the median survival time in those not cured is relatively short, so cohort effects will play only a limited role. We therefore implemented the second method in the DYNAMO-HIA model.

3.2.2 Calculating excess mortality rate from median survival

To calculate the excess mortality rate $M(t)$ from the median survival, we assume that for the highest age class in the model (95+) $M(t)$ can be calculated as:

$$M(t) = \frac{\ln(2)}{\text{medianSurvival}} \quad (7)$$

Once $M(t)$ is known for a range of ages up to the highest age, it can be found for the next lower age by solving the equation:

$$e^{-M(t+\text{int}(\text{medianSurvival}(t))*\text{mod}(\text{medianSurvival}(t)))} \prod_{s=t}^{s=t+\text{int}(\text{medianSurvival}(t))-1} e^{-M(s)} = 0.50 \quad (8)$$

where $\text{int}(T)$ is the integer part of T , and $\text{mod}(T)$ the non-integer part. Example: $\text{int}(2.34)=2$ and $\text{mod}(2.34)=0.34$.

In cases where $M(t)$ changes very fast with age, it is possible that this equation has no solution.

In that case the following error message is given: "median survival rates for age X and sex {0/1} for disease {diseaseName} are inconsistent with the median survival at the next higher age groups. Please check the data and make sure that survival does not decrease abruptly over age".

3.2.3 Calculating the probability of disease in an individual

As all input in the DYNAMO-HIA model is specified by age and gender, all calculations are carried out separately for each combination of age and gender. For clarity, however, we will leave out age and gender indexes in all the following descriptions.

The general approach here is that the prevalence odds on disease in an individual can be calculated as a relative risk times a baseline prevalence odds for that disease, making the assumption that the prevalence odds ratio of a disease in two groups is equal to the ratio of incidences rates of that disease in those groups[5, 6]. The ratio of the incidence rates are the relative risks that are given by the user. The baseline prevalence odds can be calculated by constraining the overall prevalence rate in the population (calculated by using this baseline odds for each individual and then averaging over the population), to be equal to the input prevalence rate (as given by the user). This approach is based on that used in the RIVM chronic diseases model[7]. However, in that case the assumption is that the ratio of probabilities rather than the odds is equal to the ratio of incidence ratios.

Calculation of prevalence odds for independent diseases:

Step 1: finding an initial estimate

Find an initial estimate by assuming that the prevalence relative risk in each risk factor group is equal to the incidence relative risk, and then calculate a baseline prevalence rate which is used as an initial estimate for an iterative procedure to find the baseline prevalence odds:

$$\text{Baseline odds}(d) = RR_{\text{mean}}(d) / \text{odds}(d) \quad (9)$$

To do so, we first calculate the mean relative risk in the population, and from that the baseline prevalence rate.

The mean relative risk in the population is calculated as:

For categorical risk factors:

For each disease:

$$RR_{mean} = \sum_c RR_c P(c) \quad (10)$$

where $P(c)$ is the proportion of the population having class c of the risk factor and RR_c is the relative risk on the disease for this class.

For compound riskfactor

Generate a sample of $19+N_r$ persons (20 persons in the class which has 20 duration categories (<1 year, 1-2 years, ..., 19 years and more) and N_r-1 in the other classes). Calculate the proportion of the population in each of these categories $P(\text{duration})$ using the risk factor prevalence and the distribution of the duration part of the compound risk factor as given by the user. Calculate the relative risk in each duration class:

$$RR_i(\text{duration}) = RR_{end} + (RR_{begin} - RR_{end}) \exp(-\beta * \text{duration})$$

where i is an index for each simulated person, duration is the value of the duration (0 for <1 year, 1 for 1-2 year, 19 for 19 years or more).

Then

$$RR_{mean} = \sum_{c \neq \text{durationClass}} RR_c P(c) + \sum_i P(\text{duration}) RR_i(\text{duration}) \quad (11)$$

where $P(c)$ is the proportion in class c , $P(\text{duration})$ is the probability of being in the duration class with the particular duration, RR_c the relative risk in class c and $RR_i(\text{duration})$ is the relative risk for subject i based on the duration given above.

For continuous risk factors:

Generate a sample of N_{sim} persons following the distribution of the risk factor. N_{sim} is equal to the number of simulated persons in the simulated population (as given by the user). However, in the parameter-estimation module this is increased to 100 in case a lower number was given.

To generate such a sample for a continuous risk factor, we use the mean, standard deviation and skewness of the risk factor given as input values.

For the normal distribution, mean and standard deviation define the normal distribution and skewness is zero. If skewness is non zero, and a normal distribution has be specified, the distribution is assumed to be log-normal. In that case a warning is written to the log file stating:

"normal distribution asked, but as skewness is not equal to zero, lognormal distribution is used".

In the reverse case, where a log-normal distribution is specified, but for all ages the skewness is zero, a normal distribution is used, and a the following warning is written to the log file:

"log-normal distribution asked, but as all skewness are zero, normal distribution is used"

For the lognormal distribution we have:

$$\begin{aligned}
 \text{mean} &= \exp(\mu + 0.5\sigma^2) \\
 \text{variance} &= [\exp(\sigma^2) - 1]\exp(2\mu + \sigma^2) \\
 \text{skewness} &= [\exp(\sigma^2) + 2]\sqrt{\exp(\sigma^2) - 1}
 \end{aligned}
 \tag{12-14}$$

From the skewness we can calculate the parameter σ , then from σ and the input variance we can calculate μ and from that the mean. This is the mean for a log-normal distribution starting at value 0. However, in our case we might want to use a lognormal distribution starting at a higher value (for instance, BMI will never have value 0, but will be at least 10). So the difference between the mean as given by the user and the mean calculated from the variance and skewness will be used as an offset value for the mean, that is the value at which the log-normal distribution will start.

Generating the sample of nSim persons with a continuous risk factor value R(i) will be done using the following algorithm:

1. For i=1 to Nsim generate Z(i)= (i-0.5)/Nsim
2. Calculate R(i)=F⁻¹ (Z(i)), where F⁻¹ be the inverse of the cumulative probability function of the risk factor distribution (normal or log-normal)
3. Calculate RR_i as RR_i=RR(per unit)^{(R(i) - reference value)}

And then

$$RR_{\text{mean}} = \frac{1}{N_{\text{sim}}} \sum_i RR_i
 \tag{15}$$

Step 2: calculate probability of the disease given the baseline odds

From the trial baseline odds value calculate for each class or simulated person:

$$P_c(d | c) = \frac{RR_c(d) * \text{Baseline odds}(d)}{RR_c(d) * \text{Baseline odds}(d) + 1} \quad (16)$$

or

$$P_i(d | i) = \frac{RR_i(d) * \text{Baseline odds}(d)}{RR_i(d) * \text{Baseline odds}(d) + 1} \quad (17)$$

Step 3: make overall prevalence equal to the sum of the simulated probabilities

Find the value of the baseline odds that gives (for categorical, compound and continuous risk factors, respectively):

$$\begin{aligned}
 P(d) &= \sum_c P_c(d)p(c) \\
 P(d) &= \sum_c P_c(d)p(c) + \sum_{duration} P_{duration}(d)p(duration) \\
 P(d) &= \frac{1}{Nsim} \sum_i P_i(d)
 \end{aligned} \quad (18-20)$$

using the Newton-Raphson algorithm, and where P(d) is the population prevalence of disease d as given by the user.

Calculation of baseline odds for dependent diseases:

The procedure to calculate the baseline odds for these diseases is largely the same as for independent diseases, only the relative risks for a simulated group of person are calculated differently, as they must include also the dependency on the intermediate disease.

To include this dependency, we need the prevalence of the intermediate disease as a function of risk factor levels. For these, we will use the same approximation as above, namely that for each disease the prevalence odds ratio's are equal to the incidence rate ratio's.

The odds of an intermediate diseases d_{int} can be calculated from the baseline odds of d_{int} and the relative risks.

For a class risk factor or a compound risk factor (where c here and further on will include each different length of stay in the duration class as a separate class) this is:

$$P(d_{int} | c) = Odds(d_{int}, c) / (Odds(d_{int} | c) + 1) = \frac{baselineOdds(d_{int})RR_c(d_{int})}{baselineOdds(d_{int})RR_c(d_{int}) + 1} \quad (21)$$

and for a continuous risk factor this is:

$$P(d_{int} | i) = Odds(d_{int} | i) / (Odds(d_{int} | i) + 1) = \frac{baselineOdds(d_{int})RR_i(d_{int})}{baselineOdds(d_{int})RR_i(d_{int}) + 1} \quad (22)$$

For each disease calculate a trial value of the prevalence odds by calculating RRmean as:

For categorical and compound risk factors:

$$RRmean = \sum_c RR_c(dep)P(c) \prod_{d_{int}} [P(d_{int} | c)(RR_{d_{int}}(dep) - 1) + 1] \quad (23)$$

where P(c) is the proportion of the population having class c of the risk factor and $RR_c(dep)$ is the relative risk on the particular dependent disease for this class, d_{int} are the indexes for the relevant intermediate diseases, $RR_c(d_{int})$ is the relative risk on the intermediate disease d_{int} for this class, and $RR_{d_{int}}(dep)$ is the relative risk of the intermediate disease on the dependent disease.

For continuous risk factors:

Use the sample of nSim persons as generate above.

Then calculate for each disease:

$$RR_{mean} = \frac{1}{nSim} \sum_i RR_i \prod_{d_{int}} [P(d_{int} | i)(RR_{d_{int}}(dep) - 1) + 1] \quad (24)$$

The trial version of the baseline odds can be calculated from this RRmean again using

$$Baseline\ odds(d) = RR_{mean}(d) / P(d)$$

where P(d) is the population prevalence of the disease as given by the user.

Once a baseline odds is known, we can calculate for each class or simulated person the probability on d (P(d|c) or P(d|i) depending on the type of risk factor) by averaging the probability on d given each particular combination of intermediate diseases $j_1 \dots j_n$ (thus weighting by the probability of the combination of diseases):

$$P(d | c) = \sum_{j_1=0}^{j_1=1} \sum_{j_2=0}^{j_2=1} \dots \sum_{j_n=0}^{j_n=1} \frac{RR_{dis}^* * RR_c(d) * baselineodds(d)}{RR_{dis}^* * RR_c(d) * baselineodds(d) + 1} \prod_{m=1}^{m=n} P(d_m = j_m) RR_{dis}(d | d_m = j_m) \quad (25)$$

where

$$P(d_m = 1) = \frac{RR_c(d_m) * Baseline\ odds(d_m)}{RR_c(d_m) * Baseline\ odds(d_m) + 1} \quad (26)$$

$$P(d_m = 0) = 1 - \frac{RR_c(d_m) * Baseline\ odds(d_m)}{RR_c(d_m) * Baseline\ odds(d_m) + 1} \quad (27)$$

$$RR_{dis}^* = \prod_{m=1}^{m=n} RR(d | d_m = j_m) \quad (28)$$

$RR_{dis}(d|d_i=1)$ is the relative risk on disease d for someone with disease d_i relative to someone without d_i ,

$RR_{dis}(d|d_i=0)=1$.

For a continuous risk factor the same applies, changing the index c to i .

Then the Newton-Raphson algorithm is again applied to estimate the baseline odds ratio that yields in the simulated population exactly the input population prevalence.

3.3 Estimation of the parameters of the transition rates.

In the equations in table 1, the probability terms P_i are values for an individual at a particular moment in time: these probabilities are stored during simulation as the “DYNAMO state values” of the simulated individuals. The RR(risk factor) values are calculated from the risk factor state, as given on page 11. All other terms, as well as the RR values necessary to calculate RR(risk factor) are **model parameters**. The model assumes that the model parameters are constant over time, but can differ by (attained) age and gender. The model parameters have to be supplied to the central core of the model, that uses them to calculate the transition rates (using the equations given in table 1).

Table 2 show the parameters of the transition rates. The model parameters “relative risks on diseases” and (if user-specified) “transition rates of the risk factor” are directly derived from input given by the user. The other model parameters (Baseline incidence rates, Baseline other mortality rates and Baseline fatal disease rates, the relative risks for other mortality, the attributable mortality rates, non user-specified risk factor transition rates, drift of the risk factor variation and offset, and non-user-specified risk factor drift) need to be estimated first. Below we will describe the estimation of each of those parameter.

All data are specified by age and gender. Therefore the estimation is also performed by age and gender. We will omit these in the description to simplify notation.

Table 2: Parameters of the transition rates

Type of transition rate	Parameter	Calculated from
Risk factor	Transition probability (given)	Given as input
	Transition probability (calculated)	Risk factor prevalence rates, RR for all cause mortality Baseline all cause mortality
	meanDrift, stdDrift offSetDrift (continuous risk factor)	User specified drift (optional), Risk factor mean, standard deviation and skewness, RR for all cause mortality Baseline all cause mortality
Without disease – with disease (independent diseases)	RR(riskfactor)	Given in input (categorical risk factors) or calculated from input RR information and current risk factor status. See page 11
	Baseline incidence	Prevalence of risk factor Relative risks for risk factor, see page 11
Without disease – with disease (dependent diseases)	RR(riskfactor)	Given in input (categorical risk factors) or calculated from input RR information and current risk factor status, see page 11
	Baseline incidence	Prevalence of risk factor Relative risks for risk factor (page 11) Prevalence of intermediate disease Relative risks on intermediate disease for risk factor
	RR _{id}	Given as input
To death	AM	RRs all cause mortality, all cause mortality rate, disease prevalence rates, cured fractions, fraction of acutely fatal disease, relative risks on diseases, relative risk of intermediate disease on other diseases
	RR _{oc}	Same as for AM
	Baseline other cause mortality	Same as for AM
	RR _d (riskfactor)	Given in input (categorical risk factors) or calculated from input RR information and current risk factor status.
	Baseline FD	Same as for baseline incidence plus the percentage acutely fatal disease

3.3.1 Estimation of baseline incidence rates

The baseline incidence rates is calculated as:

$$\text{baseline incidence} = \frac{\text{incidence}}{\text{RRmeanInHealthy}} \quad (29)$$

where incidence is the population incidence rate as given by the user, and RRmeanInHealthy is the average relative risk in the population without the disease.

Calculation of RRmeanInHealthy is largely equivalent to the calculation of the RRmean that is given before in section 3.2. The only difference is that incidence rates only apply to those persons who do not have the disease (as those who already have the disease are no longer at risk to get the disease), so the RR in these persons is 0.

We will use the symbols RR_c , RR_i , $P(d|c)$, and $P(d|i)$ as before. $P(d|c)$, and $P(d|i)$ are calculated as described in section 3.2.3.

For diseases with cured fraction, the healthy population are those without either the cured or the non-cured form of the disease.

Calculation of RRmeanInHealthy for independent diseases:

For categorical and compound risk factors:

For each disease:

$$RR_{mean} = \frac{1}{1 - P(d)} \sum_c (1 - P(d | c)) RR_c P(c) \quad (30)$$

where $P(c)$ is the proportion of the population having class c of the risk factor and RR_c is the relative risk on the disease for this class. For the compound risk factor, the classes c here include a separate class for each length of stay in the duration class (up to 19 years).

For continuous risk factors:

Again use nSim simulated persons from the continuous risk factor distribution

Then calculate for each disease:

$$RR_{mean} = \frac{1}{(1 - P(d))nSim} \sum_i (1 - P(d | i))RR_i \quad (31)$$

where i is an index for each simulated person

Calculation of RRmean for dependent diseases:

For categorical/ compound risk factors:

For each disease:

$$RR_{mean} = \frac{1}{1 - P(d)} \sum_c (1 - P(d | c))RR_c(dep)P(c) \prod_{d_{int}} [P(d_{int} | c)(RR_{d_{int}}(dep) - 1) + 1] \quad (32)$$

where P(c) is the proportion of the population having class c of the risk factor and RR_c(dep) is the relative risk on the particular dependent disease for this class, d_{int} are the indexes for the relevant intermediate diseases, RR_c(d_{int}) is the relative risk on the intermediate disease d_{int} for this class, RR_{d_{int}}(dep) is the relative risk of the intermediate disease on the dependent disease.

For continuous risk factors:

Calculate for each disease:

$$RR_{mean} = \frac{1}{1 - P(d)} \frac{1}{nSim} \sum_i (1 - P(d | i))RR_i \prod_{d_{int}} [P(d_{int} | i)(RR_{d_{int}}(dep) - 1) + 1] \quad (33)$$

where RR_i(dep) is the relative risk on the particular dependent disease for the simulated person, d_{int} are the indexes for the relevant intermediate diseases, RR_i(d_{int}) is the relative risk on the intermediate disease

d_{int} for this person and $RR_{d_{int}(dep)}$ is the relative risk of the intermediate disease on the dependent disease.

3.3.2 Estimation of baseline fatal disease rates

Basically the idea is to split disease incidence in incidence of acutely fatal disease and incidence of chronic disease (fatal and non-fatal incidence). Note that the incidence of chronic disease applies only to those not yet having the disease, while fatal incidence can happen also in those already having the disease. In principle fatal and non-fatal diseases can have different relative risks, but we will assume them to be equal here. This is merely because of lack of data, as from a modelling point of view separate relative risks could have been incorporated just as easily, provided that the right incidences and case-fatality rates are also available.

The input asked for by the model is the fatal fraction, the percentage of incidence (both first and recurrent events) that is fatal. We will assume in our implementation that total fatal incidence is equal to total incidence times a fatal fraction, while non-fatal incidence is equal to incidence times (1-fatality fraction), thus ignoring the differences in populations at risk.

This approximation will mean that during simulation those who already have the disease will get a slightly higher incidence of fatal disease than those without the disease, because the former will have higher risk factor levels. This is probably more realistic as an equal overall incidence of fatal disease. Other options for calculation are possible, but not warranted given the large uncertainties in the data.

After having split the overall incidence in a fatal incidence and a non-fatal incidence, the baseline non-fatal incidence rate is estimated in the same way as the baseline incidence rates for the other diseases. The fatal incidence is calculated as:

$$\text{baseline fatal incidence} = \frac{\text{fatal incidence}}{RR_{mean}} \quad (34)$$

where RR_{mean} is calculated in the same way as done for calculating the initial estimate for the prevalence odds ratio in section 3.2.3, using RR_{mean} instead of $RR_{meanInHealthy}$

3.3.3 Estimation of attributable mortality rates and other cause mortality

The input of the DYNAMO-HIA model is the excess mortality rate of each disease, as given by for instance the DISMOD program. This rate is defined as the difference of the mortality rate in those with the disease, and in those without the disease. In the case of cancers, median survival for the population that dies of the cancer can be given instead.

The model uses the method used in the RIVM CDM (Hoogenveen et al 2009) to recalculate these excess mortality rates into attributable mortality rates by removing mortality due to excess comorbidity that is already accounted for elsewhere in the model, that is:

- the comorbidity due to the fact that diseases have (modelled) risk factors in common
- the comorbidity due to the presence of a dependent disease

If diseases are independent, there is no excess comorbidity, and attributable mortality has only to be adjusted for the effects of common risk factors.

Below we will give the procedure that is used to estimate the attributable mortality and mortality due to other causes. Like elsewhere, this procedure is applied for each age and gender group separately.

Input to this procedure are:

1. excess mortality
2. all cause mortality
3. relative risks on all cause mortality for the risk factor
4. the distribution of diseases and disease combinations conditional on risk factors in the population, that can be calculated as given in section 3.2.3. Here we will use the general notation $p(d|r)$, as a general notation for both $p(d|i)$ and $p(d|c)$ as used before.

First, excess mortality (Em) from disease d is (by definition) the difference in mortality M between those with and without the disease:

$$Em(d) = M(d) - M(\text{not } d) \quad (35)$$

Given that the probability of having d in the population is equal to $p(d)$, the all cause mortality M_{tot} is equal to:

$$M_{tot} = p(d) * M(d) + (1 - p(d)) * M(\text{not } d) \quad (36)$$

Equation 35 and 36 can be combined to:

$$M(d) = M_{tot} + Em(d)(1 - P(d)) \quad (37)$$

Second, the all cause mortality and the relative risks for mortality can be used to calculate $M_{tot}(r)$, the mortality for someone with risk factor r , using the following procedure:

1. Use the input relative risks on all cause mortality to calculate a mean relative risk (RR_{mean}) in the population (using the same formulae as used for calculating RR_{mean} in section 3.2.3)
2. Use this RR_{mean} to calculate a baseline all cause mortality rate:

$$\text{baseline mortality} = \frac{\text{mortality}}{RR_{mean}} \quad (41)$$

3. Calculate $M_{tot}(r)$ by multiplying the baseline (all cause) mortality rate with the relative risk for the particular risk factor value

Now we have for each disease a value for $M(d)$, and for each risk factor class a value for $M(r)$.

In case of continuous risk factors, we here use a series of risk factor values (at least 100), generated in the same way as described at in section 3.2.2.

For both $M(r)$ and $M(d)$ we can write down the equation for this mortality in terms of the model parameters “Attributable Mortality of disease d ” ($Am(d)$), “other mortality” (dependent on risk factor status) ($Moc(r)$) and “fatal mortality” ($CF(d)$ or $CF(r)$):

As an example we write this down for 2 specific dependent diseases ($d1$ and $d2$) that both depend on intermediate disease $d3$ but not on other diseases ($d4$ to dn). This however can be generalized to more diseases, provided that there are only 2 “causal layers” (a disease that depends on another disease can not itself be a cause of another disease).

In this example, the mortality for $M(d1) = M(d1=1)$ is:

$$\begin{aligned} M(d1=1) = & am(d1) + am(d2) * P(d2|d1) + .. am(dn) * P(dn|d1) \\ & + Moc(r=1)P(d1|r=1) + \dots + Moc(r=k)P(d1|r=k) \\ & + Fatalinc1(d1) + \dots + Fatalincn(d1) \end{aligned} \quad (39)$$

With

- $P(d2|d1)$: the probability of having disease in those having disease 1
- $P(d1|r=k)$: the probability of having disease in those where risk factor $r=k$
- $Moc(r)$: the other cause mortality in those with risk factor class r
- $Fatalincn(d1)$: the average incidence of immediately fatal disease dn in those with $d1$

We can write down such a formula for each disease.

Second we can write down a similar formula for each $M(r)$:

$$M(r) = am(d1)P(d1|r) + am(d2)*P(d2|r) + .. am(dn)*P(dn|r) + Moc(r) + Fatalinc1(r) + .. + Fatalincn(r) \tag{40}$$

with

- $Fatalincn(r)$: the average incidence of immediately fatal disease dn in those with risk factor state r

This gives us a linear system of $n+c$ equations (n = number of diseases, c number of riskfactor values), containing n $am(d)$ -parameters, and c $Moc(r)$ parameters that need to be estimated. Gathering all am and Moc terms at one side of the equation, we get a linear system of equations, in which only the $am(d)$ and Moc values are unknown, and which are solved using standard linear algebra, after calculating the fatal incidence terms $Fatalincn(d1)$ and $Fatalincn(r)$ as described below.

Calculating the fatal incidence terms

For a disease $d1$ that depend on disease $d3$, $Fatalinc1$ depend on both $d3$ and r . A disease does not only depend on another disease if that other disease is a causal factor, but also in the case that they share a causal factor. This make the calculation of these terms rather complicated.

Below we give therefore general equations:

$$Fatalincj(R = r) = M_{0,CF_j} \sum_{C \in \chi} Pr(\chi = C | R = r) RR_{r \rightarrow d_j} \prod_{k: c_k=1} RR_{c_k \rightarrow d_j} \tag{41}$$

$$\text{Fatalinc}_j(d_i = 1) = M_{0,CF_j} \sum_{\mathcal{C} \in \mathcal{X}} \sum_r \Pr(\mathcal{X} = \mathbf{C}, R = r \mid d_i = 1) RR_{r \rightarrow d_i} \prod_{k:c_k=1} RR_{c_k \rightarrow d_j} \quad (42)$$

$\text{Fatalinc}_j(d_i = 1)$ Incidence of acutely fatal disease j , given that disease i is present

$\text{Fatalinc}_j(R = r)$ Incidence of acutely fatal disease i , given risk factor state r

M_{0,CF_i} Baseline fatal incidence rate of disease i .

\mathbf{C} Vector of causal disease states with elements c_k

\mathcal{X} The set of all possible vectors \mathbf{C} (all possible combinations of causal disease states)

$RR_{r \rightarrow d_i}$ The relative risk of risk factor state r on disease i

$RR_{c_j \rightarrow d_i}$ The relative risk of disease j on disease i

Given the baseline fatal incidence (estimated in section 3.3.2), these terms can be calculating by summing over the simulated population, as the prevalence of all combinations of risk factors and diseases have been estimated in section 3.2.3.

Estimation without a relative risk for mortality

In the DYNAMO-HIA model, having a relative risk for mortality is optional. Without this relative risk, mortality per risk factor class only differs because of the difference in prevalence of disease or differences in fatal mortalities. In that equation 39 is replaced by:

$$\begin{aligned} M(d1=1) = & \text{am}(d1) + \text{am}(d2) * P(d2|d1) + \dots + \text{am}(dn) * P(dn|d1) \\ & + M_{oc} \\ & + \text{Fatalinc}_1(d1) + \dots + \text{Fatalinc}_n(d1) \end{aligned} \quad (43)$$

were M_{oc} no longer depends on risk factor state, and thus is the same for every individual.

As a further equation we have:

$$\begin{aligned}
 M = & am(d1)P(d1)+am(d2)*P(d2)+..am(dn)*P(dn) & (44) \\
 & + Moc \\
 & + Fatalinc1 + .. + Fatalincn
 \end{aligned}$$

where $Fatalinc_n$ is the average fatal incidence of disease n in the population, and M the overall mortality rate in the population.

This gives us a linear system of $n+1$ equations (n = number of diseases), containing n $am(d)$ -parameters and Moc as the parameters that need to be estimated. Gathering all am and the Moc terms at one side of the equation, we get a linear system of equations, in which only the $am(d)$ and Moc values are unknown, and which are solved using standard linear algebra, after calculating the fatal incidence terms $Fatalinc_n$, by averaging the terms $Fatalinc_n(r)$ (calculated in the previous section) over the population.

Estimation in case excess mortality is zero.

If attributable mortality is zero, a small positive excess mortality will result, because those with disease have a worse risk factor profile as those without the disease (unless the disease does not depend on the risk factor, but then it would not be needed in the model).

If we reverse this reasoning, then if an excess mortality of zero is specified, a negative attributable mortality will be calculated. DYNAMO-HIA assumes that in this case the user means to request a zero attributable mortality instead of a zero excess mortality, and sets the attributable mortality to zero. In that case, for every disease for which the attributable mortality is put to zero, the system of equations will have one less equation. Estimation of the other attributable mortalities and of other cause mortality remains the same. Also in cases where the user specifies a excess mortality that is lower than the excess mortality resulting from a zero AM, the AM is set to zero.

3.3.4 Estimation of baseline other mortality and the relative risks for other mortality

In the previous section, we estimated the other mortality for each risk factor class.

For categorical risk factors, the other mortality in the reference category (where the relative risk is equal to 1), equals the baseline other mortality. The RR for other mortality then can be calculated simply as:

$$RR \text{ other mortality } (c) = \frac{Moc(c)}{\text{Baseline other mortality}} \quad (46)$$

For continuous risk factors, we fit a regression model through the estimated values of Moc(r) in order to estimate a single relative risk value on other mortality per unit of exposure (= per unit of the risk factor) as needed by the model.

Hereto we regress the log of the estimated other mortality values Moc(r) on the risk factor values R(i). The exponent of the intercept from this regression then gives the baseline other mortality, the exponent of the regression coefficients of the risk factors give the relative risks.

As this procedure forces a particular shape on the data, the average relative risk estimated by the regression, will be slightly different from the average relative risk estimated in the first part of this procedure, and because of that also the overall mortality calculated from these relative risks will not be exactly equal to the overall mortality as given by the user. Therefore the baseline other cause mortality is calibrated so that the total mortality in the simulated population becomes equal to the input total mortality. This is done by increasing or decreasing the baseline other cause mortality with the amount necessary in order to let the total mortality become equal to the total mortality specified in the input data for this age and gender.

For compound risk factors we first calculate the baseline other mortality rate and the relative risks in the same way as for the categorical risk factors, ignoring the duration information. Second we use non-linear regression to estimate the relative risks (at the begin and at the end) and the alpha parameter for the category with the duration component by fitting the following model:

$$RR_{om} = RR_{om-end} + (RR_{om-begin} - RR_{om-end}) \exp(-\alpha_{om} * \text{duration}) \quad (46)$$

As for the continuous risk factor, the average relative risk estimated for the category with the duration component by the regression will be slightly different from the relative risk estimated in the first part of this procedure, so also here the baseline other cause mortality is recalibrated so that the total mortality in the simulated population becomes equal to the input total mortality.

In this procedure we also apply some restrictions on RRend and RRbeing: When RR-end for total mortality (as given by the user) is equal to the RR of total mortality in one of the other classes (without a

duration component), RR-end for other mortality is restricted to be equal to RR-other mortality for this class(es).

Similarly, when RR-begin for total mortality is equal to the RR for total mortality of one of the other classes (without a duration component) RR-begin other mortality is restricted to be equal to the RR for other mortality for this class. In case both RR-end and RR-begin for total mortality are equal to an RR from other classes, both are restricted to be equal to those RRs, and only alpha is fitted.

In some cases, however, such a restriction is not possible. For example, when the RR-begin is fixed to 0.85, and the RR-end to 1, while the RRs in the duration category on average are below 0.85 or above 1. In that case the restriction on RR-begin is removed. In the case that regression does not return an estimate even after lifting this restriction, also the restriction on the RR-end is removed. In case these restrictions are lifted for these reasons, messages are written to the log-file.

Another possible complication is that the user puts all persons in the starting population in a single duration class. In that case no estimation of time dependence is possible, and RR-end is made equal to RR-begin, implying no time dependence. In that case a warning message is given:

"100% of the initial population has the same duration. Therefore it is not possible to estimate a time dependent other mortality or other disability. In case other mortality/disability is requested, those relative risks will be made constant over time."

A possible complication for all types of risk factors is that other mortality(i) could become negative in case of high $am(d)$ or $Fatalinc(d)$ values in combination with simulated subjects with high $P(d|i)$ values or $RR_d(i)$ values. If this only occurs sporadically this is handled by setting other mortality(i) to zero for categorical risk factors, and ignoring this class when fitting the other mortality in case of a continuous risk factor or a time in the risk factor class. The baseline other mortality then also is recalibrated, assuring that the total other cause mortality in the population remains valid.

However, if negative other cause mortality values occur frequently, this indicates that the combination of excess mortality rates and disease prevalence rates entered by the user leads to more mortality in the population than is physically possible based on all cause mortality rates.

The following warning will be given if negative mortality occurs in more than 30% of the simulated cases of the same age and sex: "Warning: Disease related mortality is too large relative to all cause mortality".

3.3.5 Estimation of Non user-specified risk factor transition rates (nett transition rates)

3.3.5.1 For a categorical risk factor

The program can estimate 1 year risk factor transition probabilities from the risk factor prevalence rates. This option is called “nett transition rates”. The procedure used is described in more detail in van de Kasstele et al.[8]. The general idea behind this estimation is that it estimates transition probabilities that will keep the age specific prevalence rate of the risk factor stable over time. That is, in future years the age distribution of the risk factor is assumed to be the same as the current distribution by age. In that case we can say that if prevalence rates in the next higher age group differs from prevalence rates in the current age group, people must have switched from one class to another. It is however impossible to estimate the transition probabilities back and forth from cross sectional data, we can only estimate the nett transition probabilities (with 2 risk factor classes A and B, the nett effect is the difference between the flow from A to B and that from B to A, thus the minimum number of persons that have to change risk factor level in order to obtain the new prevalence rate).

The estimation of transitions for such cases can be formulated as a transportation problem, a topic solved in operations research. A transportation problem basically attempts to find best way to fulfil the demands of a number of demand points using the supplies from a number of supply points, under the consideration of the costs involved in shipping the product from a supply point to a demand point.

A transportation problem is specified by the supply, demand, shipping costs, a decision variable, and an objective function. In this case the supplies are the predicted class prevalence rates $\tilde{\mu}_0$ based on the prevalence rates in the current year. It is a vector of length n . The demands are the class prevalence rates $\tilde{\mu}_1$, also a vector of length n , in the next year. The costs \mathbf{C} , an $n \times n$ matrix, are the costs involved in moving from one class to another. The decision variable \mathbf{x} , also an $n \times n$ matrix, contains the transitions between classes from row i to column j , and is the variable of interest. The objective function J is the total shipping cost from classes in the current year to classes in the next year. This function has to be minimized:

$$\min J = \sum \mathbf{C} * \mathbf{x}, \quad (47)$$

subject to

$$\sum_{j=1}^n \mathbf{x}_{\cdot,j} = \tilde{\boldsymbol{\mu}}_0, \sum_{i=1}^n \mathbf{x}_{i,\cdot} = \tilde{\boldsymbol{\mu}}_1, \text{ and } x_{ij} \geq 0, \text{ for } i, j = 1 \dots n. \quad (48)$$

The above constraints entail that the row sums of \mathbf{x} , i.e. people from a specific classes, must be equal to the supply $\tilde{\boldsymbol{\mu}}_0$. The column sums of \mathbf{x} , i.e. people to specific classes, must be equal to the demand $\tilde{\boldsymbol{\mu}}_1$. This also means that the total supply must be equal to the total demand. The last constraint entails that no transition may be negative.

Some assumptions have to be made about the cost matrix \mathbf{C} , since no direct “costs” can be assigned to people moving from one class to another. It is assumed that most people tend to stay in their class. Some will switch one class up or down, some will switch two or even more classes, however we assume this becomes very unlikely. We therefore assign zero costs to transitions to the same class, a cost of one unit for one class up or down, a cost of three units for two classes up or down, and so forth. This choice of costs will assure that it is more likely that two persons change one class, than one person changing two classes, as the latter costs 3 units, and the first 2 times 1 unit.

The transportation problem can now be easily solved by a linear programming algorithm, for which we use the simplex method. Once the transitions \mathbf{x} have been found, the transition probability matrix \mathbf{M} is then given by

$$\mathbf{M} = (\mathbf{x} / \tilde{\boldsymbol{\mu}}_0)^T, \quad (49)$$

so that a new state is given by $\tilde{\boldsymbol{\mu}}_1 = \mathbf{M} \tilde{\boldsymbol{\mu}}_0$. A transition probability matrix is calculated for each age.

In the method given above selective mortality is included by not using $\tilde{\boldsymbol{\mu}}_0$, the prevalence of the current age, in the algorithm, but $\tilde{\boldsymbol{\mu}}_0^*$, the hypothetical prevalence after one year, taking mortality into account, but not letting persons change state:

$$\tilde{\boldsymbol{\mu}}_0^* = \frac{(\exp(-\mathbf{RR} * \text{baselineMort})) \tilde{\boldsymbol{\mu}}_0}{\sum_j \exp(-\mathbf{RR}_j * \text{baselineMort}) \tilde{\mu}_{0j}} \quad (50)$$

For \mathbf{RR}_j the user specified relative risks are used.

For a compound risk factor

For a compound risk factor the same method is used as given for a classified variable. However, here we first need the RR for all cause mortality for the class with duration. This RR is calculated as:

$$RR_{duration} = \sum_d RR_d P(d) \quad (51)$$

where $P(d)$ is the probability of being in duration class d , given that one is in the class with duration.

In principle the duration information itself could also be used to calculate net transition probabilities. However, information on duration is often much more difficult to obtain, and expert opinion might be used instead of real data. Therefore this road was not taken.

For a continuous risk factor

For a continuous risk factor a similar method is used. For a continuous risk factor we do not have separate states, therefore it is more intuitive to look at it from the point of the update rule applied to an individual.

The update rule for a **normally distributed** variable with mean μ_0 at the current age and μ_1 at the next age and matching standard deviations σ_0 and σ_1 is:

$$\begin{aligned} X(1) &= X(0) + (\mu_1 - \mu_0) + \text{stdDrift} * \varepsilon & \text{if } \sigma_1 > \max(\sigma_0) \\ X(1) &= X(0) + (\mu_1 - \mu_0) & \text{if } \sigma_1 \leq \max(\sigma_0) \end{aligned} \quad (52)$$

Where ε is a randomly generated draw from the standardized Normal distribution, and the parameters of this update rule are the “drift” in mean ($\mu_1 - \mu_0$) and stdDrift is a constant with which the randomly drawn value of ε is multiplied. The latter is put to zero if the standard deviation at a particular age is equal or smaller to that on the previous age. The effect of this way to update the continuous risk factor is that in case the standard deviation of the risk factor distribution is increasing over age, this increase is reproduced by the random component. Note that this procedure assumes that the standard deviation has been smoothed over age. Using unsmoothed figures will result in stdDrift values that are too large.

The drift in the mean causes all individuals to increase their risk factor value by the same amount: that means there is, based on only the mean drift (and without the random component), perfect tracking over age.

Values of the mean drift are calculated directly from the mean of the risk factor by age:

The mean drift can be calculated as the difference of the risk factor mean at the particular age and the next age. Like with the nett transitions for the categorical risk factor, we first adjust the mean at the age for selective mortality by calculating μ_0^* , the mean of the distribution of the risk factor x after applying the risk factor specific mortality rate for one year:

:

$$\tilde{\mu}_0^* = \frac{\int x \exp(-RR(x)baselineMort)P(x)dx}{\int \exp(-RR(x)baselineMort)P(x)dx} \quad (53)$$

Similarly, we also calculate the standard deviation σ_0^* , the standard deviation of the distribution of risk factor x after applying the risk factor specific mortality rate for one year:

$$(\tilde{\sigma}_0^*)^2 = \frac{\int (x - \tilde{\mu}_0^*)^2 \exp(-RR(x)baselineMort)P(x)dx}{\int \exp(-RR(x)baselineMort)P(x)dx} \quad (54)$$

The stdDrift then in calculated as:

$$stdDrift = \sqrt{\sigma_1^2 - (\sigma_0^*)^2} \quad (55)$$

The standard error drift is made the same for all scenario's, the user can only change the amount of the meandrift.

The numerical integration procedure used in Dynamo uses 100 equidistant intervals between $\mu_0 - 4\sigma_0$ to $\mu_0 + 4\sigma_0$. Preliminary simulations showed that this range yields the most accurate results given the use of 100 intervals, and is superior to using equal probability intervals.

For a **lognormally** distributed risk factor, the risk factor value is first logtransformed into a normally distributed variable. The update rule on the transformed variable is identical to that given above for a normally distributed risk factor. After applying this part of the update rule it then is back transformed. Drift and change in variance on the transformed can be calculated from the input data in exactly the same

way as for normally distributed variables, but after logtransformation. If there is an offset, a possible “offset drift” is added after back-transformation. The offset drift can be calculated as the difference of the offsets of two consecutive years.

Both the offset drift and the standard error drift are the same for all scenario’s, the user can only change the amount of the meandrift.

The meandrift should always be given by the user as the absolute change (increase/decrease) in the mean value, also in the case that a log-normal distribution is specified for the data. In case of a log-normal distribution, the program calculates the amount of change needed on the logscale.

This way to parametrize the update rule means that during the first step of simulation, the risk factor value changes in the way the user will expect, that is, it increases with the amount given. In case of the log-normal distribution, the model uses the relative change of the amount (above the offset) to update the risk factor. So when the mean value of the risk factor at a particular age changes during simulation (for instance due to extra selective mortality in the scenario), the absolute drift will change, as a relative change is applied by the model. Therefore, using the lognormal distribution could lead to behaviour that is not always intuitive for the user.

Selective mortality is taken into account in the same way as for the normally distributed risk factors, that is by calculating the parameters of the distribution of the risk factor on the log scale after taking one year survival into account. For this calculation we assume that the offset is not affected by mortality.

4 Generation of Initial Population

4.1 Population characteristics

Each person in the simulation population has the following characteristics:

- A** age
- G** gender
- C** value $\{1,2,\dots,n_r\}$ of classified risk factor: present if risk factor type = classified or compound
- R** value of continuous risk factor, or the length in the duration class: present if risk factor type = continuous or compound
- DM₁ to DM_{nc}** Disease probability tables for each cluster of diseases (1 to n_c). For non-cancers each table has size $2^{(N_in_Cluster(c))}$. Each value will be indicated by $P(i,j)(j \leq 1)$, giving the probability of the disease combination conditional on being alive. Therefore the sum of the values of a single table is always 1, and therefore only $2^{(N_in_Cluster(c))-1}$ values need to be stored.
For diseases with a cured fraction (non-lethal cancers) the table has size 2. For independent diseases (not related to other diseases) the table has size 1.
- S** Probability of survival.

Input needed for the generation a an initial population

	Input needed
Input on structure	Number of disease clusters, number of diseases in each cluster, whether a cluster concerns a disease with a cured fraction, type of risk factor (compound, continuous, categorical) nSim: numbers in the initial population per age and gender stratum, minimal and maximum age of individuals at the start of simulation
Input Parameters (directly given by the user)	Relative risks on diseases $RR(r)$, $RRdis(di,dd)$
Intermediate Parameters	Baseline prevalence odds(d) (see 3.2.3),Cured

	prevalence fraction (d) (see 3.2.1)
--	-------------------------------------

The initial population is generated by:

1. constructing a distribution of risk factors that is close to that of the population being simulated.
2. calculating the probability of disease in each simulated individual based on the risk factor state of this person.

At the start of the simulation all persons are alive, so the survival probability is set to 100%.

Below we will describe the algorithms used for the generation of the initial population in detail.

4.2 Generation of characteristics for the reference population

The reference population is the population to which no intervention is applied. For this population the characteristics are generated as described below.

Generation of A and G:

For each combination of A and G in the run, we create Nsim persons.

Generation of C:

We will generate a number of subjects in each risk factor class close to the real prevalence rates using the following algorithm:

1. For each category: Calculate the number of subjects in this category from nSim, rounded downwards:
 $N_c = \text{Int}(P(c) * n\text{Sim})$, but make this minimally one person
2. Calculate the number of missing subjects in the simulation: $n_todo = n\text{Sim} - \text{Sum}(N_c)$
3. Calculate for each c: $P_todo_c = (P(c) * n\text{Sim} - \text{Int}(P(c) * n\text{Sim})) / n_todo$; This is the proportion of “todo” subjects that should go to class c;
4. Randomly draw the n_todo persons using P_todo_c.

Generation of R:

Generate a sample of Nsim persons following the distribution of the risk factor using the following algorithms:

For continuous risk factors:

1. For $i=1$ to Nsim generate $Z(i) = (i-0.5)/Nsim$
2. Calculate $R(i) = F^{-1}(Z(i))$, where F^{-1} be the inverse of the cumulative probability function of the risk factor distribution

For compound risk factors:

Generate C as given above for categorical risk factors, but include the 20 duration classes as part of P(c). Make $R = 0$ for all classes without duration, and $R=0$ to 19 for the duration class.

Generation of DM

Once we have simulated a risk factor value for each person in the simulation, we use this value to generate the initial probability of disease. Below we give the detailed algorithms used for this

Generation of DM for a single independent disease:

1. Calculate RR_i as :
 - $RR_i = RR(\text{per unit})^{(R(i) - \text{reference value})}$ (continuous)
 - $RR_i = RR(c)$ (categories or non-duration categories)
 - $RR_i = RR_{\text{end}} + (RR_{\text{begin}} - RR_{\text{end}}) \exp(-\beta * R(i))$ (duration categories)
2. Calculate Odds(d, i) as $\text{BaselineOdds}(d) * RR_i$,
3. Calculate $P(d,i)$ as $\text{Odds}(d,i) / (\text{Odds}(d,i) + 1)$

Generation of DM for diseases with a cured fraction:

1. Check whether this cancer is independent of other diseases
2. If not, ask the user to change this
3. Calculate RR_i as

$$RR_i = RR(\text{per unit})^{(R(i) - \text{reference value})} \text{ (continuous)}$$

$$RR_i = RR(c) \text{ (categories or non-duration categories)}$$

$$RR_i = RR_{\text{end}} + (RR_{\text{begin}} - RR_{\text{end}}) \exp(-\beta * R(i)) \text{ (duration categories)}$$

4. Calculate $P(d\text{-total}, i)$ as above
5. Calculate $P(d\text{-notcured})$ as $P(d\text{-total}, i) * (1 - \text{CuredPrevalenceFraction}(d))$
6. Calculate $P(d\text{-cured}, i)$ as $P(d\text{-total}, i) * \text{CuredPrevalenceFraction}(d)$

Generation of DM for a cluster of dependent diseases:

1. With n diseases in the cluster ($\text{clustersize} = n$), there will be $2^n - 1$ characteristics in the simulation to characterize the disease state of the cluster. The combination: no diseases is $1 - \text{sum}(\text{all other combinations})$, so it is not explicitly stored.
2. Calculate RR_i for each d as
 - $RR_i(d) = RR(d, \text{per unit})^{(R(i) - \text{reference value})}$ (continuous)
 - $RR_i(d) = RR(c, d)$ (categories or non-duration categories)
 - $RR_i(d) = RR_{\text{end}}(d) + (RR_{\text{begin}}(d) - RR_{\text{end}}(d)) \exp(-\beta * R(i))$ (duration categories)
3. Loop through all combinations of diseases in the cluster and calculate the probability of each combination, which is :

$$P(\text{combi}) = \prod_{\text{dep}} P(\text{dependent disease} \mid \text{all independent diseases}) \prod_{\text{indep}} P(\text{independent disease})$$

where $P(\text{(in)dependent disease})$ is the probability that the disease has the value it has in the combination. $P(\text{dependent disease} \mid \text{all independent diseases})$ and $P(\text{independent disease})$ can both be calculated in with the following equation, given that $RR_{\text{dis}}(d_1=1 \mid d_2=1)$ (the relative risk on disease 1 given disease 2) is always 1 in cases of independent diseases:

$$P(d=1 \mid \text{other diseases in cluster}) =$$

$$\begin{aligned}
 \text{For } d = 1: & \frac{\text{BaselineOdds}(d) * RR_i(d) \prod_{d1 \in cluster} (d1 = 1 \text{ in combi}) RRdis(d = 1 | d1 = 1)}{1 + \text{BaselineOdds}(d) * RR_i(d) \prod_{d1 \in cluster} (d1 = 1 \text{ in combi}) RRdis(d = 1 | d1 = 1)} \\
 \text{For } d = 0: & 1 - \frac{\text{BaselineOdds}(d) * RR_i(d) \prod_{d1 \in cluster} (d1 = 1 \text{ in combi}) RRdis(d = 1 | d1 = 1)}{1 + \text{BaselineOdds}(d) * RR_i(d) \prod_{d1 \in cluster} (d1 = 1 \text{ in combi}) RRdis(d = 1 | d1 = 1)}
 \end{aligned}$$

Generation of S:

Value S will be 1 for every person in the initial population

4.3 Generation of characteristics for scenario populations

Dynamo offers three possibilities for scenarios:

1. Using the same initial population as for the reference population, but changing the transition probabilities for the risk factor
2. Using the same transition probabilities for the risk factor, but changing the risk factor distribution in the initial population
3. Changing both the risk factor distribution in the initial population, and the the transition probabilities for the risk factor

In the first situation, the initial population of the reference scenario is simply copied.

In the second situation, for continuous risk factors, the model assumes that the ranking of persons within the population remains the same. So the individual with the xth rank in the old distribution (that is, the distribution in the reference population) gets the value of the xth rank in the scenario distribution.

For categorical and compound risk factors, subjects need to change risk factor status. In this case, we calculate the minimum number of persons that need to change risk factor status in order to reach the new risk factor distribution, using the same algorithm as is used to calculate net transition rates (see section 3.3.5).

For each non-zero risk class transition, we generate scenario subjects by applying each transition to every individual in the “old” risk factor class to which that transition applies. An example might clarify. Say there are 100 individuals in the reference simulation population, and the old prevalence is 50% in class 1, 40% in class 2 and 10% in class 3, and the new prevalence is 20% in class 1, 20% in class 2 and 60% in class 3. Then we assume that 20 individuals originally in class 1 have to go to class 2, 10 of those originally in class 1 have to go class 3 and all 20 subjects in class 2 have to go to class 3. So the non zero transitions are : $1 \rightarrow 2$, $1 \rightarrow 3$ and $2 \rightarrow 3$. In that case the simulation population exists of:

- 50 individuals originally in class 1, with risk factor changed to 2
- the same 50 individuals originally in class 1, with risk factor changed to 3
- 40 individuals originally in class 2, with risk factor changed to 3.

During post-processing these individuals receive weights of 20/50, 10/50 and 40/40 respectively, also adding unchanged individuals (from the reference population) from classes 1, 2 and 3 with weights 20/50, 0/50 and 10/10.

This procedure makes that for these scenario’s stable results can be obtained with a relatively low number of simulated persons, especially in the case that only a moderately number of persons changes risk class, which is often the case in Health Impact Assessments.

In the third case, individuals from the reference population can not be reused in this way as part of the scenario population, as also transition rates differ between reference and alternative scenario. In this case, for a categorical or compound risk factor, the transitions from old to new prevalence are used to randomly draw a new value for every individual in the reference population. This procedure means that the new prevalence is not fixed, but affected by a random component. In these scenario’s one needs to make N_{sim} sufficiently large.

For continuous risk factors, a new value is assigned in the same way as before, retaining the ranking of subjects in the population.

In all cases, all scenario-individuals are directly matched to a particular individual in the reference population. Matched individuals will share the same random number during simulation.

The initial disease state of the individual does not change by changing the exposure, so that the scenario-individual at the start of simulation has the same disease state as the matched reference individual. The only exception are the newborns (age 0 at the start of simulation), that get immediately the prevalence belonging to the new risk factor state. In practise, disease prevalence is low at birth, so this choice will hardly influence the outcomes.

5 Description of simulation module

During the simulation, there will be simulated persons with the following characteristics

A	age
G	Gender
C	value {1,2,...,n _r } of classified risk factor: present if risk factor type = classified or compound
R	value of continuous risk factor or the length of stay in the duration class for compound risk factors
DM₁ to DM_{n_c}	Disease probability tables for each cluster of diseases (1 to n _c). For non-cancers each table has size $2^{(N_in_Cluster(c))-1}$. Each value will be indicated by P(s), giving the probability of the disease combination conditional on being alive. For diseases with a cured fraction the table has size 2. For independent diseases (not related to other diseases) the table has size 1.
S	Probability of survival

Update rules are used to update these values in steps of one years. Section 5.1.3 describes these update rules.

In the DYNAMO model, simulation starts with a population aged 0-95 (unless the user restricts the age to a smaller range of simulation). All simulated persons in the starting population are simulated until at the age of 105. This makes it possible to calculate cohort life-expectancies for this cohort. (see chapter 6). Simulated persons born during simulation are only simulated for the number of time steps needed to reach the last simulation year that was given by the user.

5.1.1 Synchronising of scenarios

It is essential for a an efficiently working program that all scenarios that need to be compared are run with the same random numbers for the same draws (“parallel universe approach”)[9].

Therefore random draws need to be synchronized between scenarios. This is done by matching each individual in the simulation of an alternative scenario to an individual of the reference scenario, and then starting these matched individuals with the same random seed in all scenarios. We referred to this as the

“parallel universe” approach, as we simulate what would happen to an individual in a parallel universe, where all random happenings are identical to the real universe, but another policy has been being implemented. In making the initial populations for the scenario’s, individuals are copied from the reference population to a scenario population, including their random seed.

In the implementation the use of the random-generator is not conditional on characteristics of an individual, so that random number generation can not get out of sync between matched reference and scenario individuals.

5.1.2 Adding newborns to the population

Generating newborns is handled by including individuals with negative age in the initial population. With the exception of age itself, which is increased by one year during each time step, states of the newborns are kept constant until they reach age 0. From that moment on they are considered as born, and risk factor, disease and survival state are updated using the update rules described below.

5.1.3 Update rules for characteristics of the simulated persons.

5.1.3.1 Update rule for A

Input:

<u>Parameters</u>	<u>Characteristics</u>
None	A

Rule:

$$A \text{ updated} = A + 1$$

5.1.3.2 Update rule for G

Input:

<u>Parameters</u>	<u>Characteristics</u>
	G

Rule:

G updated = G

5.1.3.3 Update rule for C

Input:

<u>Parameters</u>	<u>Characteristics</u>
Transition matrix for age A, gender G	C (A and G needed to select parameter but not in rule itself)

Rule:

Draw a random C from the vector of probability Q. Q is found by taking the row of the transition matrix(A,G) belonging to current risk factor value C.

When zero transition are specified, C stays constant.

5.1.3.4 Update rule for R (if type=continuous risk factor)

Input:

<u>Parameters</u>	<u>Characteristics</u>
meandrift (A,G) , stddrift(ratio)(A,G), offsetdrift(A,G) , type (normal or lognormal), offset(A,G)	R (A and G needed to select parameter but not in rule itself)

Rule:

If R has a normal distribution

R updated=R +drift(A,G) +Normal*stddrift(A,G)

Where Normal is a draw from the standardized normal distribution

If R has a log-normal distribution:

R updated=offset+offsetdrift+exp(log(R-offset) +drift(A,G) + Normal*stddrift(A,G))

NB: for the lognormally distributed R, meandrift(A,G) is not the drift of the mean as entered into the program, but the drift of the parameter of the lognormal distribution.

5.1.3.5 Update rule for R (if type = compound risk factor)

Input:

<u>Parameters</u>	<u>Characteristics</u>
Number of class of which this is the duration	R ,C

Rule:

If C= duration class then R updated = R+1

Else R=0

Update rule for DM(d) (for independent disease)

This rule needs to be repeated for each disease d

Input:

<u>Parameters</u>	<u>Characteristics</u>
Baseline incidence: $i_0(d,A,G)$, Attributable mortality $AM(d,A,G)$, Relative risks for risk factor on the disease (RR),	DM(d), R, C (A en G needed to select parameter but not in rule itself)

Rule:

Calculate RR as:

$RR = RR(\text{per unit})^{(R(i) - \text{reference value})}$ (in case of a continuous riskfactor)

$RR = RR(c)$ (in case of a categorical risk factor, or the non-duration classes of a compound risk factor)

$RR = RR_{\text{end}} + (RR_{\text{begin}} - RR_{\text{end}}) \exp(-\beta * R(i))$ (in case of the duration category of a compound risk factor)

Calculate the incidence Inc as

$$Inc = i_0(d, A, G) * RR$$

Solving the exponential of the transition rate matrix

$$\begin{bmatrix} -Inc & 0 \\ +Inc & -AM(d, A, G) \end{bmatrix}$$

yields the following formulae for DM(d)

DM(d) updated =

$$\frac{(DM(d)AM(d, A, G) - Inc)e^{(Inc-AM(d, A, G))} + Inc(1 - DM(d))}{(DM(d)AM(d, A, G) - Inc)e^{(Inc-AM(d, A, G))} + AM(d, A, G)(1 - DM(d))}$$

In case AM(d,A,G)=I(d,A,G) the denominator becomes zero. In this case we solve the system of differential equations for the transition matrix

$$\begin{bmatrix} -Inc & 0 \\ +Inc & -Inc \end{bmatrix}$$

which yields:

$$1 - \frac{1 - (DM(d))}{1 + Inc(1 - DM(d))}$$

Update rule for DM(d) (for diseases with cured fraction)

Input:

<u>Parameters</u>	<u>Characteristics</u>
Baseline incidence: $i_0(d, A, G)$, Attributable mortality AM(d,A,G), Relative risks for risk factor on the disease (RR), percentage cured (CuF(d,A,G))	DM(d), R, C (A and G needed to select parameter but not in rule itself)

DM has two values:

DM(c): the probability of having cured disease

DM(nc): the probability of having a lethal form of the disease.

In terms of a transition matrix, there are 3 states: without the cancer, with cured cancer, with non-cured cancer. The vector D_{old} of initial probabilities (before the update) of these 3 states are:

1-DM(c)-DM(nc); DM(c),DM(nc).

The transition matrix T is (first row gives the rate of change of the first state as a function of the states) :

$$\begin{bmatrix} -Inc_{tot} & 0 & 0 \\ +Inc_{tot} f_{cured} & 0 & 0 \\ +Inc_{tot} (1 - f_{cured}) & 0 & -AM(d, A, G) \end{bmatrix}$$

Here inc_{tot} is the incidence of the disease, calculated in the same way as in the update rule for a single independent disease. The matrix exponential of this matrix also has an analytical solution, which is implement in the software.

The general analytic solution for DM(c) (cured prevalence) after a timestep t is:

$$\left\{ \frac{\left(e^{at} (a - itot) (fcured (-1 + ptot) + e^{itott} (fcured + pcured - fcuredptot)) \right)}{e^{itott} itot} \left(-1 + fcured + pcured - fcuredptot - e^{at} (fcured + pcured - fcuredptot) \right) + a \left(e^{at} (-1 + fcured) (-1 + ptot) + e^{itott} (-pcured + ptot) + e^{(a+itott)t} (fcured + pcured - fcuredptot) \right) \right\}$$

and for DM(nc) (non cured prevalence):

$$\left\{ \frac{\left(e^{at} (-1 + fcured) itot (-1 + ptot) + e^{itott} (a (-pcured + ptot) + itot (-1 + fcured + pcured - fcuredptot)) \right)}{e^{itott} itot \left(-1 + fcured + pcured - fcuredptot - e^{at} (fcured + pcured - fcuredptot) \right) + a \left(e^{at} (-1 + fcured) (-1 + ptot) + e^{itott} (-pcured + ptot) + e^{(a+itott)t} (fcured + pcured - fcuredptot) \right)} \right\}$$

where $pcured=DM(c)$ at the start of the time step;

ptot = the sum of DM(c) and DM(nc) at the start of the time step and

a= attributable mortality. In the DYNAMO model with implemented this solution for t=1.

The incidence itot is calculated in the same way as for independent diseases.

5.1.3.6 Update rule for DM(s) (for clusters of mutually dependent diseases)

Input:

<u>Parameters</u>	<u>Characteristics</u>
Baseline incidence: i_0 (d,A,G), Attributable mortality AM(d,A,G), Relative risks for risk factor on the disease (RR),Relative risks for diseases on diseases (RRdis)	DM(s), R, C (A en G needed to select parameter but not in rule itself)

DM(s) gives the disease probabilities (conditional on being alive) for each possible combination of the diseases (state s) in the cluster. We use s as argument here in stead of d to indicate that this characteristic indicates a state, and not a single disease.

The state: “none of the diseases” is not stored in the simulation as it can be calculated as 1- sum(other entries). Thus for n diseases in the cluster there are $2^n - 1$ values of DM(s), where s is a number, having values 1 to 2^n , which written in binary form indicates the presence of each disease. For instance, s= 100100¹ indicates that the first and fourth disease are present, but diseases 2,3,5 and 6 are not present.

For this rule there are the following steps:

1. add state s=0 and calculate DM(0) as 1 minus the sum of all DM(s).
2. calculate the **rate of going** from state s1 to each other state s2 $T(s1 \rightarrow s2)$. (see later for how to calculate these rates), and construct the matrix **T** that contains these values.
3. calculate the matrix of **probabilities of going** from state s1 to each other state s2 **P** (containing all probabilities $P(s1 \rightarrow s2)$) from

¹ In the software implementation the ordering is the other way round: the last digit represents disease number 0, the one before that disease 1 etc.

$$\mathbf{P} = \exp(\mathbf{T}.t)$$

This is carried out by applying the algorithm described by Gallivan et al.[4].

4. calculate the probability of being in state s_2 after update, conditional on being alive at the **beginning** of the interval:

$$DMU(s_2) = \sum_s DM(s)P(s \rightarrow s_2)$$

5. calculate the probability of being in state s_2 after update, conditional on being alive at the **end** of the interval :

$$DM(s_2)_{updated} = DMU(s_2) / \sum_s DMU(s)$$

Calculation of $T(s_1 \rightarrow s_2)$

To calculate the matrix with transition rates, $T(s_1 \rightarrow s_2)$ we first determine which bits (diseases) in s_2 differ from those in s_1 , as those are the diseases that need to change during the transition. For this we apply the following rules:

1. The model does not include remission, so in case there is a disease present in s_1 that is not present in s_2 , $T(s_1 \rightarrow s_2) = 0$
- 2 For the other transition rate between two different states s_1 and s_2 : Let O_{ac} be the set of diseases which are 0 in s_1 , but 1 in s_2 (the set of diseases that is to be acquired) and O_{-1} the set of diseases which are 1 in s_1 (irrespective of their value in s_2). As the rate gives the change in an infinitely small period of time, there can only be one change in disease between state. So if O_{ac} has more than one member, $T(s_1 \rightarrow s_2) = 0$.

If O_{ac} has one member, d , then let $inc_0(d)$ be the incidence of disease d , calculated in the same way as of the update of a single independent disease (including the dependence on risk factor values, but not yet the dependence on the presence on other diseases). The rate of acquiring disease d during the update is given by:

$$T(s_1 \rightarrow s_2) = inc_0(d) \prod_{e \in O_{-1}} RR_{dis}(e, d)$$

3 Mortality rates (attributable mortality and case fatality) are entered as (negative) terms to the diagonal of the matrix. In case a disease in the combination is zero, also the incidence to this disease is a negative contribution to the diagonal term. So the transition rate $T(s1 \rightarrow s1)$ for a diagonal term is:

$$T(s1 \rightarrow s1) = -\sum_s T(s1 \rightarrow s) - \sum_d I(d=1)am_d - \sum_d fatal_d(s1)$$

where am_d is the attributable mortality for disease d (added for those disease that are 1 in $s1$) and $fatal_d(s1)$ the fatal incidence due to disease d in those with state $s1$ (added for all states). This fatal incidence is calculated in the same way as given for the incidence above.

5.1.3.7 Update rule for S

Input:

<u>Parameters</u>	<u>Characteristics</u>
Baseline incidence: $i_0(d,A,G)$, Attributable mortality $AM(d,A,G)$, Relative risks for risk factor on the disease (RR), Relative risks for diseases on diseases (RRdis) Relative risks for risk factor on other cause mortality (RRom) Baseline other cause mortality (baselineOM)	R, C, D, all DM (A en G needed to select parameter but not in rule itself)

Rule:

Total survival is given by the multiplication of survival fractions due to the different clusters of diseases, and due to other causes.

For independent diseases the survival fraction due to the disease (that is the fraction with which the disease decreases survival in a period t) is given by:

$$S_d = \frac{am(1 - DM(d)) e^{-inc} + (am \cdot DM(d) - inc) e^{-am}}{am - inc}$$

With:

Am: the attributable mortality for the disease

Inc: the incidence of the disease, calculated as given in the update rule for a single independent disease
 DM(d): the prevalence of the disease (=DM(d))

The survival fraction for other cause mortality is given by:

$$S_{OM} = e^{-om}$$

where om is other cause mortality. This is calculated as follows:

First calculate RR_{om} as :

$$RR_{om} = RR_{om} (\text{per unit})^{(R(i) - \text{reference value})} \text{ (in case of a continuous riskfactor)}$$

$$RR_{om} = RR_{om} (c) \text{ (in case of a categorical risk factor, or the non-duration classes of a compound risk factor)}$$

$$RR_{om} = RR_{om-end} + (RR_{om-begin} - RR_{om-end}) \exp(-\beta * R(i)) \text{ (in case of the duration category of a compound risk factor)}$$

Then $om = RR_{om} \cdot \text{Baseline-OM}$.

For diseases with a cured fraction the survival fraction due to these diseases is calculated from the values DM(cured) and DM(not cured) at the start of the time step as:

$$S_d = DM(cured) + e^{-am} \cdot DM(notcured) + (1 - DM(cured) - DM(notcured)) \cdot \left\{ e^{-inc_cured} + (1 - e^{-inc_cured}) \frac{(inc_cured - inc_noncured)}{inc_cured} + \frac{(e^{-inc_cured} - e^{-am}) inc_noncured}{am - inc_cured} \right\}$$

For clusters of diseases, the survival is calculated as:

$$S_c = \sum_s DMU(s)$$

where DMU is calculated as given in the update rule for clusters of diseases.

Finally S is calculated by multiplying all the survival fractions as calculated above.

6 Post-processing module and output generated

The DYNAMO simulation module delivers a population of simulated persons. As output in this form can be very large, it is not stored or made available to the user. This output contains the actually simulated persons, and is not always representative of the population that is simulated. In the post-processing stage the simulated persons are upweighted to the real population numbers as given by the user in the input module. After this weighting, detailed group specific data can be outputted in the form of excel readable XML files. Also DYNAMO-HIA can be run in batch-mode, in which case a csv file is generated, containing similar information, but in one flat file rather than in a workbook with multiple worksheets. In batch mode also output objects are created and stored, which can be viewed later using the graphical user interface. In the batch mode, only the data generated over the user specified running time are included, despite the fact that the model simulates all subjects alive at the start of simulation until age 105.

6.1 Weighting procedure

Weighting is applied to upweight the simulated population to numbers in the real population. First, for categorical risk factors, weights are applied to make the simulated population representative in terms of risk factor classes and changes between scenarios in risk factor classes (see section 4.3). Due to this weighting process, after weighting the risk factor prevalence of categorical risk factors at the start of simulation will be exactly equal to the prevalence entered into the model, with the exception of scenarios populations where both the initial prevalence of the risk factor and the risk factor transitions are changed. The latter will contain a random component.

Second, the number per age and gender are up-weighted to become representative of the population, using the demographic data as supplied by the user.

This procedure delivers population data for scenario populations where the intervention is 100% successful and is applied to every age and gender group. From these data, effects of interventions with lower success rates or effects of an intervention reaching only particular age or sex groups are calculated by combining “successful” scenario individuals and “unsuccessful” reference individuals in the correct proportion.

6.2 Content of DYNAMO-output files

excel readable XML files can be read by using “open with” and then choose “excel”.

The output in these files are in the form of numbers of persons in the total population.

The excel readable output exists in two forms:

Yearly data:

For each simulated year (as given by the user), the numbers in the simulated population (males, females or both combined) by:

- Age
- Actual risk class in that year (for continuous variables: categorized in groups)

For each age /risk factor class combination, the files contain:

1. For continuous risk factor: the mean value
2. For compound risk factor: the mean value of duration
3. The total number of years lived in that class
4. Either numbers in the class with the different disease states (= combinations of diseases within each cluster), or the numbers with each disease (to choose by the user).
5. The number of DALY-weighted years in that class
6. The number of persons with one or more diseases in that class

Each year is on its own worksheet of the excel workbook.

Cohort data:

For each simulated cohort (group with same age at start of simulation), the numbers by
Year of simulation (always up to the age of 105)

Risk class at the start of follow-up (for continuous variables: categorized in groups)

For each age/gender/risk factor class combination, the files contain:

1. For continuous risk factor: mean value
2. For compound risk factor: mean duration
3. Total number of persons alive

4. Disease state or disease (to choose by the user).
5. The number of DALY-weighted years in that class
6. The number of persons with one or more diseases in that class

Each cohort is on its own worksheet of the excel workbook. The user can produce separate workbooks for men and women, or a single workbook in which numbers from men and women are combined.

Batch output:

For each simulated year (as given by the user), the numbers in the simulated population (males and females separately) by:

- Age
- Actual risk class in that year (for continuous variables: categorized in groups)

For each age /risk factor class combination, the files contain:

1. The year
2. The scenario
3. The gender
4. The risk factor class
5. For continuous risk factor: the mean value
6. For compound risk factor: the mean value of duration
7. The total number of years lived in that class
8. Either numbers in the class with the different disease states (= combinations of diseases within each cluster)
9. The number of DALY-weighted years in that class
10. The number of persons with one or more diseases in that class

Correspondence with the use of data in the graphical user interface

Yearly data are in the DYNAMO-HIA graphical user interface - output module used to calculate:

- Prevalence of disease by year, or by age (within each year)
- Risk factors distribution by year, or by age (within each year)
- Mortality numbers (those dying between the year and the next year) by year, or by age (within each year)
- Period life expectancy and Sullivan Health expectancy

Cohort data are in the DYNAMO-HIA graphical user interface - output module used to calculate:

- Cohort Life expectancy (by age)
- Disease free life expectancy and life expectancy with disease (cohort)
- Survival
- Disability adjusted life expectancy (cohort)

6.3 Calculation of integral measures of health.

The program calculates the following health expectancy values:

- total life expectancy
- life expectancy free of a particular disease
- life expectancy free of all diseases in the model
- disability free life expectancy or DALE. Which of these is calculated, depends on the data given by the user. If DALY-weights were given, DALEs are calculated. If percentages of disability have been given, disability free life expectancy is calculated.
- Number of life years gained or lost compared to the reference scenario
- Number of life without disease gained or lost compared to the reference scenario
- Disability weighted life years gained or lost compared to the reference scenario

For life expectancy free of all modelled diseases, the prevalence of having no disease is calculated for each simulated individual. As the clusters of diseases that are calculated in DYNAMO are mutually independent given the risk factor history, this can be calculated by multiplying the probabilities of being disease-free from each cluster of diseases.

Similarly, within each cluster, the DALY-weights can be applied to the different disease-states in this cluster (including the healthy state, that is the state without any of the disease in the model). If a person has two diseases, the DALY weight given to the person is calculated as:

$$1 - W_{\text{daly.total}} = (1 - W_{\text{daly-disease1}}) \cdot (1 - W_{\text{daly-disease2}}) \cdot (1 - W_{\text{daly-healthy}})$$

For ease of explanation, we will use the term “ability” for 1-DALY-weight, defined as either 1 – the fraction disabled or as a health valuation where 1 is perfect health, and 0 health as bad as being dead. Then in general the ability due to different diseases is calculate as:

$$ability = baseline\ ability \prod_d ability_d$$

where the baseline ability is ability in a person without any of the diseases that are included in the model, and $ability_d = 1$ for a person without disease d . The abilities of the different health states in the cluster can be averaged (weighted for the probability of the health state) to the average “ability” in the cluster. As disease clusters are mutually independent given risk factor history, for a simulated individual the overall ability can then be calculated by multiplying the abilities due to different disease clusters.

Health expectancies are calculated both as cohort health expectancies and cross-sectional (“Sullivan”) health expectancies.

The cohort health expectancies are calculated directly from the simulating the cohort (always simulated until age 105), and the calculation consists of simply counting the years lived in the particular health state during this follow-up period by all simulated individuals, and dividing this by the number of individuals at the start of simulation. Implicitly this means that the life expectancy in the age category 105+ is taken as 0.5. However, as individuals of age 105 and older are extremely rare, this choice does not influence the results.

Cross-sectional health expectancy is a synthetic measure, that is calculated for a particular calendar year. Thereto first the mortality in that year is calculated by taking the difference between the survival in the next year and the current year. This mortality is combined with the prevalence of disease prevalence or disability in the particular year in a period life table constructed based on the data from the one calendar year, and a health expectancy calculated from the life table in the usual way. For the cross-sectional health expectancy, only ages upto 95 are included. The life expectancy in the age category 95+ is taken as: $-1/\ln(1-\text{mortality rate}(95))$.

From the cohort life expectancy also the number of life years lost in the population under an alternative scenario compared to the reference scenario are calculated in two ways. First, the number of life years lost or gained can be calculated for those have a particular age at the start of simulation (say all 10 year olds). This is done by simply multiplying the cohort life expectancy at age 10 with the number of 10 year olds in the population. Second, this can be done for the entire population by adding up the years gained or lost from all ages.

Similarly, years without disease gained or lost and disability weighted years gained or lost are calculated by multiplying the cohort life expectancy without disease and the disability weighted life expectancy respectively with the number of persons in the population at the start of simulation.

The disability weighted years gained or lost for a single age are in magnitude comparable to the DALY measure: only DYNAMO calculates this prospectively (years lost over a prospective life), while the classical DALY is a constructed cross-sectional measure. Also, unlike the classical DALY, the DYNAMO “DALY” takes competing morbidity into account.

Reference List

1. Lhachimi, S.K., et al., *Standard tool for quantification in health impact assessment a review*. Am J Prev Med, 2010. **38**(1): p. 78-84.
2. Boshuizen, H.C., et al., *The DYNAMO-HIA model: An efficient implementation of a risk factor / chronic disease Markov model for use in health impact assessment.(submitted)*
3. Moler, C. and C. Van Loan, *Nineteen Dubious Ways to Compute the Exponential of a Matrix*. SIAM Review, 1978. **20**(4): p. 801-836.
4. Gallivan, S., et al., *A computational algorithm associated with patient progress modelling*. Computational Management Science, 2007. **4**(3): p. 283-299.
5. Neogi, T. and Y. Zhang, *Re: "Easy SAS calculations for risk or prevalence ratios and differences"*. Am J Epidemiol, 2006. **163**(12): p. 1157; author reply 1159-61.
6. Alho, J.M., *On prevalence, incidence, and duration in general stable populations*. Biometrics, 1992. **48**(2): p. 587-92.
7. Hoogenveen, R.T., P.H. van Baal, and H.C. Boshuizen, *Chronic disease projections in heterogeneous ageing populations: approximating multi-state models of joint distributions by modelling marginal distributions*. Math Med Biol, 2009.
8. Kasstele J, et al., *Estimating net transition probabilities from cross-sectional data with application to risk factors in chronic disease modelling*.
9. Shechter, S.M., et al., *Increasing the efficiency of Monte Carlo cohort simulations with variance reduction techniques*. Med Decis Making, 2006. **26**(5): p. 550-3.



© RIVM 2010

Parts of this publication may be reproduced, provided acknowledgement is given to the 'National Institute for Public Health and the Environment', along with the title and year of publication.