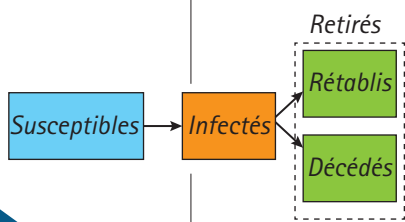


Les modèles compartimentaux

Les modèles compartimentaux, souvent utilisés en biologie mathématique sont des simplifications, parfois drastiques, de la réalité. De telles simplifications permettent d'avoir une vision d'ensemble, de capter l'essence des phénomènes et de voir la forêt et non l'arbre qui la cache.

Regardons deux exemples.

Christiane Rousseau
Université de Montréal



Le premier exemple est le modèle SIR¹ de propagation des maladies infectieuses. Il prend les hypothèses les plus simples possibles et son objectif est de décrire l'évolution d'une épidémie en fonction du temps et en ignorant les phénomènes spatiaux. Le premier ingrédient est la division de la population en trois compartiments : les personnes susceptibles, les personnes infectées et les personnes retirées (rétablies ou décédées).

Le second ingrédient est donné par l'ensemble de règles qui permettent de décrire le nombre d'individus dans chaque compartiment au jour n , étant donné un nombre initial d'infections au jour 0. Soit $S(n)$, $I(n)$, $R(n)$, les nombres de personnes susceptibles, infectées et retirées au jour n .

Règle 1 : Le nombre de susceptibles diminue du nombre de nouvelles infections.

Règle 2 : Le nombre d'infections augmente du nombre de nouvelles infections et diminue du nombre de nouveaux retraits.

Comment suivre l'évolution des quantités $S(n)$, $I(n)$, $R(n)$? Pour générer de nouvelles infections, il faut des individus contagieux. De plus, il faut absolument que des individus contagieux soient en contact avec des individus susceptibles. En effet, s'il n'y a personne à contaminer, il n'y aura pas de nouvelles infections. La manière la plus

simple de mettre cela en équation est que le nombre de nouvelles infections au jour $n+1$ soit proportionnel à la fois à $S(n)$ et à $I(n)$, donc au produit $S(n)I(n)$ avec un facteur de proportionnalité p . Ce facteur p peut être interprété comme la probabilité qu'une personne susceptible rencontre une personne infectée et que cette rencontre résulte en une infection.

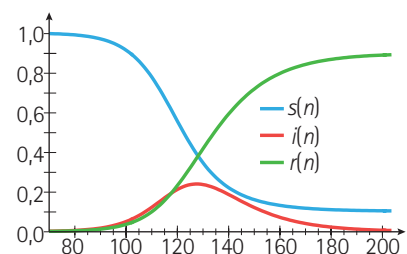
Quant à la règle 2, l'hypothèse la plus simple est qu'une fraction q des infectés est retirée chaque jour. Ceci donne :

$$\begin{aligned} S(n+1) &= S(n) - pS(n)I(n), \\ I(n+1) &= I(n) + pS(n)I(n) - qI(n), \\ R(n+1) &= R(n) + qI(n), \end{aligned}$$

que l'on peut simuler pour différentes valeurs de p et q .

Au travers des simulations se dessinent trois grandes lois :

- la *croissance exponentielle* des infections au début de la pandémie;
- le *pic des infections*, ou sommet de la vague;
- le phénomène de l'*immunité de groupe* par lequel l'épidémie s'éteint avant que toutes les personnes susceptibles n'aient attrapé la maladie.



$s(n) = S(n)/N$, $i(n) = I(n)/N$, $r(n) = R(n)/N$ où N est la population totale.

1. Voir aussi « Naviguer au travers d'une épidémie », Accromath 15.2, été-automne 2020

À cause du terme $pS(n)I(n)$, ce modèle est non linéaire. C'est ce terme qui contrôle le pic des infections et le phénomène de l'immunité de groupe qui ne pourraient se produire dans un modèle linéaire. Le modèle SIR décrit l'évolution d'une épidémie et on ne peut faire plus simple.

« Tout doit être aussi simple que possible, mais pas plus simple. », Albert Einstein, 1933

Expliquons maintenant la mécanique de ces trois lois. Pour cela, il est préférable de regarder les variables

$$s(n) = \frac{S(n)}{N}, i(n) = \frac{I(n)}{N}, r(n) = \frac{R(n)}{N}$$

qui représentent les fractions de la population de personnes susceptibles, infectées et retirées et qui prennent leurs valeurs dans $[0,1]$. Les équations deviennent :

$$\begin{aligned} s(n+1) &= s(n) - pNs(n)i(n), \\ i(n+1) &= i(n) + pNs(n)i(n) - qi(n), \\ r(n+1) &= r(n) + qi(n). \end{aligned}$$

La deuxième équation donne

$$i(n+1) = i(n)(1 + pNs(n) - q)$$

Au début de l'épidémie $s(n)$ est proche de 1. Donc, $pNs(n) \approx pN$. Par suite, en début d'épidémie,

$$i(n) \approx i(0)(1 + pN - q)^n. \quad (*)$$

Comme le n est en exposant, c'est la croissance exponentielle annoncée.

On peut aussi réécrire la deuxième équation comme

$$i(n+1) - i(n) = i(n)(pNs(n) - q).$$

Alors, on voit que $i(n+1) - i(n)$ a le signe de $(pNs(n) - q)$. Mais, au fur et à mesure que l'épidémie progresse, $s(n)$ décroît. Quand $s(n)$ est suffisamment petit, $(pNs(n) - q)$ devient négatif et les infections se mettent à décroître, expliquant le pic de l'épidémie et l'immunité de groupe.

Aucun des ingrédients du modèle SIR n'est exact. Pourtant, ce modèle est très riche d'instructions sur l'évolution d'une épidémie.

« Tous les modèles sont faux, mais certains sont utiles.² »
George E.P. Box,
statisticien, 1976

Le modèle SIR décrit ci-dessus ne prédit qu'une seule vague. Pourtant, on en a observé plusieurs pour la COVID-19. Et on a vu des individus attraper la COVID-19 plus d'une fois.

En fait, des raffinements naturels du modèle nous permettent de suivre tous ces phénomènes. Ainsi, le paramètre p représentant la probabilité qu'une personne susceptible rencontre une personne infectée et que cette rencontre résulte en une infection varie au cours du temps. Il diminue lors de mesures de distanciation sociale. Il

2. Traduction libre de « All models are wrong, but some are useful. »



augmente lorsqu'un nouveau variant, plus contagieux, apparaît. Donc, il est naturel de le regarder comme dépendant du temps, c'est-à-dire comme une fonction $p(n)$ du jour n . Ce qui est important c'est toujours l'équation

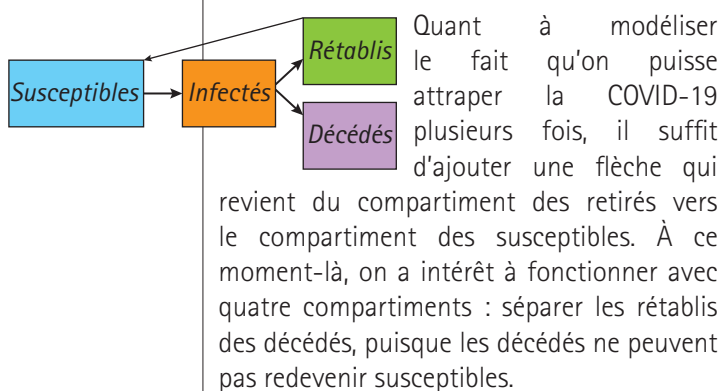
$$i(n+1) - i(n) = i(n)(p(n)Ns(n) - q)$$

qui indique comment évoluent les infections. Ainsi,

$$i(n+1) - i(n) > 0 \text{ si } p(n)Ns(n) - q > 0,$$

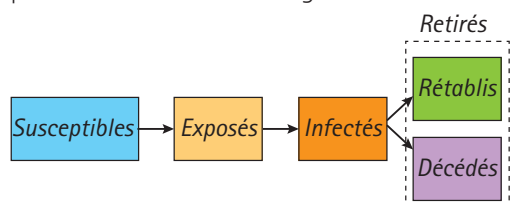
$$i(n+1) - i(n) < 0 \text{ si } p(n)Ns(n) - q < 0.$$

Et des oscillations dans la valeur de $p(n)$ peuvent permettre plusieurs vagues.

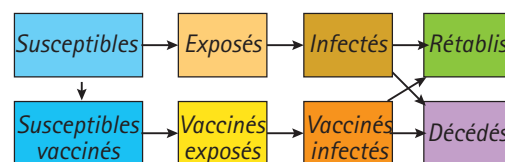


D'autres raffinements du modèle

Un premier défaut du modèle est qu'une personne devient contagieuse dès qu'elle a été infectée, alors qu'on sait qu'il y a une période de latence avant que la personne devienne contagieuse. À cause de cela le modèle prédit le pic de l'épidémie plus tôt que ce qu'il se produit en pratique. La correction utilisée par les épidémiologistes est d'introduire un compartiment de personnes exposées. Une personne infectée passe quelque temps dans ce nouveau compartiment des exposées avant de migrer vers le compartiment des personnes infectées contagieuses.



Le tournant dans la pandémie de COVID-19 qui a permis l'allègement des mesures sanitaires a été la vaccination. On peut modéliser la vaccination en ajoutant un compartiment pour les personnes susceptibles vaccinées. Pour la COVID-19, on a vu que les personnes vaccinées peuvent quand même être infectées. Il faudrait alors diviser les personnes infectieuses en deux compartiments : infectieuses vaccinées et infectieuses non vaccinées. Les paramètres p et q peuvent être différents pour les deux compartiments.



À vous d'ajouter d'autres raffinements.

Donc, on peut voir le processus de modélisation comme un processus en plusieurs étapes en partant d'un squelette qui montre la structure globale qu'on habille ensuite de couches de détails.

La modélisation mathématique permet d'avoir une vision d'ensemble

Plusieurs questions se posent. Comment déterminer la valeur numérique des paramètres? Et quelle est la fiabilité des prédictions d'un modèle mathématique pour simuler l'avenir et faire des prédictions.

Pour déterminer la valeur numérique des paramètres, il faut des données sur les nombres d'infections, la durée des infections, etc. Si la période d'infection est en moyenne de L jours, il est naturel de penser que, chaque jour, une fraction $1/L$ des personnes infectées sont retirées, et donc que le paramètre q est de l'ordre de $1/L$. Au début de l'épidémie on a vu la croissance exponentielle qui, dans la variable s'écrit

$$I(n) \approx I(0) (1 + pN - q)^n.$$

Donc, si on connaît $I(n)$ pour plusieurs n , on peut espérer estimer $1 + pN - q$, et donc p une fois qu'on a estimé q . Une difficulté est que toutes les infections ne sont pas rapportées. Et dès qu'on introduit des mesures sanitaires, cela change le paramètre p , ce qui a été le cas avec la COVID-19 mais qui était rarement le cas avec la grippe dans le passé.

Si l'on a accès à des données sur une période suffisamment longue, on simule le modèle pour cette période passée. On compare les prédictions du modèle avec les données observées pendant cet intervalle de temps. Lorsque le modèle dépend de paramètres, la technique de validation croisée consiste à diviser les données en deux ensembles disjoints : le premier ensemble est utilisé pour estimer les paramètres, et le second pour valider le modèle avec ces valeurs de paramètres estimées.

Le taux de reproduction de base R_0

On a beaucoup entendu parler pendant la pandémie du fameux R_0 , le *taux de reproduction de base*, qui correspond au nombre moyen d'infections secondaires générées par une infection primaire au début de la pandémie. Prenons un individu infecté au début de la pandémie. La journée n , il génère $p S(n)$ nouvelles infections. On a mentionné qu'il est assez naturel que q soit de l'ordre de $1/L$, où L est le nombre moyen de jours pendant lequel un individu est contagieux. Donc, en L jours, cet individu va générer $p S(n) L \approx p S(n)/q$ nouvelles infections. Mais, au début de l'épidémie, $S(n)$ est de l'ordre de N . Ceci suggère la formule,

$$R_0 = \frac{pN}{q},$$

et on voit bien dans (*) que le facteur $1 + pN - q$ est plus grand que 1 si et seulement si $R_0 > 1$.

Ce taux de reproduction de base est un grand concept unificateur. Il a été emprunté dans d'autres contextes, dont la modélisation des espèces invasives. Voyons-le dans notre deuxième exemple.

Modéliser une espèce invasive, la matricaire inodore

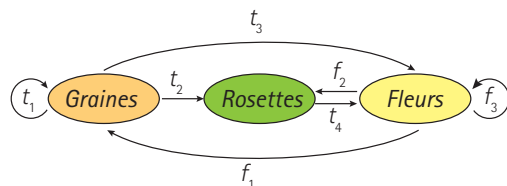
La matricaire inodore ou camomille sauvage a été introduite d'Europe. Elle se répand principalement dans les sols perturbés, ce qui lui permet d'envahir les cultures, dont elle diminue le rendement. Cette plante est un fléau dans les cultures de l'Alberta. Le conseil de recherches de l'Alberta a donc commandité une recherche en 2004-05 pour comprendre comment la plante se répand, et tenter d'optimiser des stratégies de contrôle. Cette recherche a été menée par Tomas de Camino Beck, étudiant de doctorat à l'époque, et par son directeur de recherche, le biologiste mathématicien Mark Lewis. Regardons-la.

La plante est annuelle ou bisannuelle et son cycle de vie comprend trois états : l'état de graine, l'état de rosette (ou groupe de feuilles), et l'état de fleur. La germination a lieu en mai-juin, et les graines qui ne germent pas peuvent attendre jusqu'à 15 ans dans la terre un moment favorable pour germer.

On regarde l'état de la plante en juillet et la transition d'un mois de juillet au mois de juillet suivant. L'unité de temps est donc l'année.

- Une graine peut germer l'année suivante et donner une plante à l'état de rosette ou de fleur.
- Une plante à l'état de rosette sera à l'état de fleur l'année suivante.
- Une plante à l'état de fleur va faire des graines. Trois états pour ces graines l'été suivant :
 - des graines qui dorment dans la terre,
 - des plantes à l'état de rosette,
 - des plantes à l'état de fleur.

Le schéma du cycle de vie est du même type que les diagrammes à compartiments d'une maladie infectieuse vus précédemment.



Le modèle proposé est linéaire et représenté par une matrice. Si G_n , R_n et F_n sont les quantités de graines, de plantes à l'état rosette et de plantes à l'état fleur en l'année n , alors

$$\begin{pmatrix} G_{n+1} \\ R_{n+1} \\ F_{n+1} \end{pmatrix} = \begin{pmatrix} t_1 & 0 & f_1 \\ t_2 & 0 & f_2 \\ t_3 & t_4 & f_3 \end{pmatrix} \begin{pmatrix} G_n \\ R_n \\ F_n \end{pmatrix} \quad (**)$$

Tous les coefficients de la matrice sont positifs ou nuls. Les coefficients t_i sont des coefficients de transition, donc inférieurs à 1, pour tenir compte de la mortalité et du fait que la transition est partielle. Les coefficients f_i sont des coefficients de fécondité. Comme chaque plante peut produire jusqu'à 256 000 graines, les coefficients f_i sont très grands. L'équation (**) est de la forme $P_{n+1} = AP_n$, où A est la matrice et P_n le vecteur des « populations ».

Le taux de reproduction de base R_0 pour la matricaire inodore

Le taux de reproduction de base compte le nombre moyen de « descendants directs ». Que signifie *descendant direct* dans ce contexte? Une graine peut rester graine pendant quelques années avant de germer. Donc, il faut regarder plus loin que l'année suivante. Aussi, pour parler de descendance, on veut qu'il y ait eu une forme de *fécondation*. On décompose donc la matrice A comme $A = T + F$, où T est la matrice de transition, et F , la matrice de fécondité :

$$T = \begin{pmatrix} t_1 & 0 & 0 \\ t_2 & 0 & 0 \\ t_3 & t_4 & 0 \end{pmatrix} \text{ et } F = \begin{pmatrix} 0 & 0 & f_1 \\ 0 & 0 & f_2 \\ 0 & 0 & f_3 \end{pmatrix}.$$

Les descendants directs après un an sont donnés par FP_0 . Après deux ans, ce sont ceux qui ont survécu après un an, puis se sont reproduits, donc donnés par FTP_0 . Après trois ans, ce sont ceux qui ont survécu deux ans, puis se sont reproduits, donc donnés par FT^2P_0 . Par suite, si I est la matrice identité, le nombre de descendants directs est donné par

$$F(I + T + T^2 + \dots) P_0.$$

Si T était un nombre plus petit que 1, on pourrait déduire de la formule

$$(1 + T + T^2 + \dots + T^n)(1 - T) = 1 - T^{n+1}$$

qu'à la limite

$$(1 + T + T^2 + \dots)(1 - T) = 1$$

En suivant la même logique et en remplaçant 1 par la matrice identité I , on admettra qu'une formule similaire est valide pour la matrice T dont tous les coefficients sont dans $[0, 1[$:

$$(I + T + T^2 + \dots + T^n)(I - T) = I$$

d'où l'on tire

$$(I + T + T^2 + \dots) = (I - T)^{-1}$$

où $(I - T)^{-1}$ est l'inverse de la matrice $I - T$. La nouvelle matrice $B = F(I - T)^{-1}$ donne le nombre de *descendants directs*. Il existe un nombre R_0 , qui exprime la « tendance » de la matrice B (voir encadré). Ce nombre R_0 est appelé *taux de reproduction de base*.

Ce taux de reproduction de base joue bien le rôle qu'on veut qu'il joue. La population croît si $R_0 > 1$, décroît si $R_0 < 1$ et est stationnaire si $R_0 = 1$.

De plus, R_0 est facile à calculer (voir encadré) et vaut

$$R_0 = \frac{f_1(t_3 + t_2 t_4)}{1 - t_1} + f_2 t_4 + f_3.$$

Les coefficients de la matrice A doivent être évalués sur le terrain. Ceci a été fait par Tomas de Camino Beck. Les coefficients de fécondité varient beaucoup d'une année à l'autre car la fécondité dépend des conditions de chaleur et d'humidité. Voici les résultats pour les deux années de l'étude :

$$A = \begin{pmatrix} 0,08 & 0 & 36\,376 \\ 0,27 & 0 & 517 \\ 0,04 & 0,45 & 298 \end{pmatrix} \text{ en 2004,}$$

Le taux de reproduction de base de la matricaire inodore

Le calcul de R_0 est facile, mais demande des concepts avancés d'algèbre linéaire. La matrice B a une forme bien spéciale. En effet, les deux premières colonnes de F sont nulles et sa troisième colonne est de la forme

$$\begin{pmatrix} f_1 \\ f_2 \\ f_3 \end{pmatrix}.$$

Donc, si r_3 est la troisième rangée de $(I - T)^{-1}$, les trois rangées du produit $F(I - T)^{-1}$ sont de la forme $f_1 r_3, f_2 r_3, f_3 r_3$. Un calcul donne que

$$r_3 = \begin{pmatrix} \frac{t_3 + t_2 t_4}{1 - t_1} & t_4 & 1 \end{pmatrix}.$$

$$A = \begin{pmatrix} 0,08 & 0 & 1775 \\ 0,27 & 0 & 25,24 \\ 0,04 & 0,45 & 14,53 \end{pmatrix} \text{ en 2005.}$$

Ceci donne $R_0 \approx 0,175f_1 + 0,45f_2 + f_3$.

Pour 2004, $R_0 \approx 6896$, car $0,175f_1 \approx 6366$, $0,45f_2 \approx 232$ et $f_3 \approx 298$.

Pour 2005, $R_0 \approx 335,93$, car $0,175f_1 \approx 310$, $0,45f_2 \approx 11,4$ et $f_3 \approx 14,53$.

Pour ces deux années, on voit que R_0 est très grand et que le terme $0,175f_1$ est beaucoup plus grand que les deux autres.

Utiliser R_0 pour tenter de contrôler la matricaire inodore

La très grande valeur du R_0 traduit le fait que la matricaire inodore est une plante très envahissante dont il est difficile de se débarrasser. La solution la plus simple serait de l'empêcher de faire des graines par un contrôle mécanique en enlevant les têtes de fleurs avant que les graines ne soient mûres, mais c'est une solution difficile à mettre en pratique quand la plante envahit des champs de blé.

Pour comprendre l'action de la matrice B , on change de base. De par la forme de la matrice, il existe deux vecteurs linéairement indépendants V_1 et V_2 tels que $BV_1 = 0$ et $BV_2 = 0$. Pour compléter une base, il faut un troisième vecteur V_3 . On peut le choisir tel que BV_3 soit un multiple μV_3 de V_3 .

Ce coefficient μ est le R_0 cherché¹. Il vaut

$$R_0 = \frac{f_1(t_3 + t_2 t_4)}{1 - t_1} + f_2 t_4 + f_3.$$

1. Comme les trois rangées sont multiples l'une de l'autre, la matrice a deux valeurs propres nulles. Et comme la somme des valeurs propres est égale à la trace de la matrice (c'est-à-dire la somme des coefficients de la diagonale), alors la troisième valeur propre, soit R_0 , est égale à la trace.

Aussi, le terme $0,175f_1$ est la plus importante contribution du R_0 . C'est parce qu'il y a énormément de graines qui dorment dans la terre et attendent une autre année avant de germer. Ceci enseigne donc que toute stratégie de contrôle de la matricaire inodore doit s'échelonner sur plusieurs années et qu'une stratégie limitée à une année donnée serait inefficace. Les chercheurs considèrent des stratégies de contrôle mixtes : contrôle chimique à l'aide d'herbicides, contrôle biologique avec des insectes s'attaquant aux graines ou aux plants de matricaire inodore et contrôle mécanique lorsque c'est possible.

Conclusion

Un grand nombre de phénomènes évoluent dans le temps et la modélisation mathématique permet de comprendre leur dynamique. Parmi les modèles utilisés, les modèles compartimentaux forment une large classe de modèles relativement faciles à étudier. Le taux de reproduction de base est un grand concept unificateur pour plusieurs de ces modèles. Il aide à la compréhension et à la prise de décision.