

Apprentissage Automatique
Master Sciences Cognitives

Alexis Nasr

27 mars 2020

1 Chapitre 1

2 Réseaux récurrents

3 Un réseau convolutionnel permet de prendre en entrée une séquence d'ob-
4 servations tout en l'analysant par petit bouts à l'aide d'une *fenêtre glissante*.

5 Mais un tel réseau ne permet pas de prendre en entrée une séquence de
6 longueur quelconque.

7 Nous nous sommes servi de ce type de réseaux pour des séquences de
8 caractères de longueur fixe (100, 1000, ...) dans le cas de l'identification de
9 la langue d'un texte ou d'images de 28×28 pixels pour la reconnaissance de
10 chiffres manuscrits.

11 On pourrait imaginer un réseau dont le but est de prédire, par exemple,
12 quel événement va avoir lieu à un moment donné dans un film.

13 Il n'est pas évident de concevoir un tel réseau car on ne sait pas à quels
14 moments dans le passé ont lieu des événements qui aideront à prédire un
15 événement au temps présent.

16 Afin d'autoriser la prise en compte d'événements arbitrairement distants,
17 on introduit la notion d'**état** du réseau au temps t ¹, noté \mathbf{s}_t (s pour state en
18 anglais).

19 Au temps t , le réseau a accès à l'entrée au temps t : \mathbf{x}_t , ainsi qu'à l'état
20 du réseau au temps $t - 1$: \mathbf{s}_{t-1} .

21 \mathbf{s}_i véhicule de l'information utile pour des prédictions au temps $t > i$, il
22 se présente sous la forme d'un vecteur.

23 L'état du réseau constitue en même temps une sortie et une entrée du
24 réseau : il est produit par ce dernier (sortie) et pris en compte (entrée) pour
25 des prédictions futures.

1. De manière plus général, ce qu'on appelle temps ici est en fait une position dans une séquence, tel que la position d'une lettre dans une séquence de lettres.

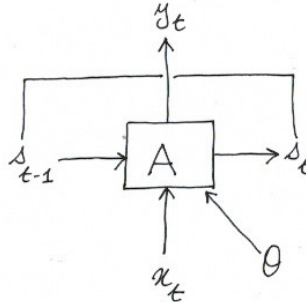


FIGURE 1.1 – un réseau récurrent

1 C'est cette caractéristique qui lui confère le nom de **réseau récurrent**,
 2 noté RNN dans la suite de ce document.

3 Un RNN est généralement représenté à l'aide d'une boucle, comme dans
 4 la figure 1.1 qui illustre bien son caractère récurrent.

5 Si on observe le fonctionnement d'un réseau dans le temps, on obtient
 6 une vision "dépliée" de ce réseau, comme dans la figure 1.2.

7 Dans la version dépliée, chaque occurrence du réseau A correspond à
 8 un instant, on voit ainsi le déroulement du réseau dans le temps : à chaque
 9 instant t (dans la figure t prend respectivement les valeurs 1, 2, 3 et 4) l'entrée
 10 \mathbf{x}_t est lue, ainsi que le vecteur d'état \mathbf{s}_{t-1} , une sortie $\hat{\mathbf{y}}_t$ est calculée, ainsi que
 11 le vecteur d'état \mathbf{s}_t .

12 Il est important de bien remarquer que, comme dans un réseau convolu-
 13 tionnel, le réseau A est "dupliqué" plusieurs fois, il s'agit du même réseau,
 14 et donc des mêmes paramètres.

15 Nous avons représenté de manière explicite les paramètres θ dans la fi-
 16 gure 1.2 pour insister sur le fait qu'il s'agit d'un même jeu de paramètres,
 17 pour toutes les étapes du traitement.

18 **Notations :**

19 — L'entrée du réseau est constitué d'une séquence $\mathbf{x}_{1..n}$ de vecteurs de
 20 dimension d_{in} .

— La sortie du réseau au temps t est un vecteur $\hat{\mathbf{y}}_t$ de dimension d_{out} :

$$\hat{\mathbf{y}}_t = rnn(\mathbf{x}_{1..n})$$

21 — Pour chaque sous-séquence $\mathbf{x}_{1..t}$, le réseau produit une sortie \mathbf{y}_t .

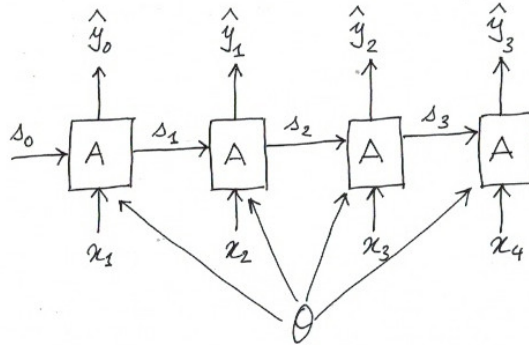


FIGURE 1.2 – un réseau récurrent déplié

1 On note rnn^* la fonction qui produit la séquence $\mathbf{y}_{1\dots t}$

$$\hat{\mathbf{y}}_{1\dots t} = rnn^*(\mathbf{x}_{1\dots t})$$

2 Le vecteur de sortie $\hat{\mathbf{y}}_t$ est généralement utilisé pour effectuer une prédiction,
3 par exemple à l'aide d'un MLP.

4 De manière plus précise, le RNN peut être défini à l'aide d'une fonction
5 R qui prend en entrée un état \mathbf{s}_{t-1} et une entrée x_t pour produire un nouvel
6 état \mathbf{s}_t :

$$\mathbf{s}_t = R(\mathbf{s}_{t-1}, \mathbf{x}_t)$$

7 On définit de plus une fonction O qui produit une sortie \mathbf{y}_i à partir de
8 \mathbf{s}_i :

$$\hat{\mathbf{y}}_t = O(\mathbf{s}_t)$$

9 Pour amorcer la pompe, on a besoin d'un état \mathbf{s}_0 que l'on considèrera
10 pour l'instant comme un vecteur nul (toutes ses composantes valent 0).

11 Le modèle présenté ci-dessus, à l'aide des deux fonctions R et O est un
12 modèle abstrait. Pour définir un réseau concret, il faut spécifier la forme
13 exacte de R et de O , ce que l'on fera plus tard.

14 La définition de R et de O confère au réseau un caractère récursif, qui
15 permet de prendre en compte, lors de la production de y_i , toutes les entrées
16 $x_{j \leq i}$.

1 Illustration pour $i = 4$

$$\begin{aligned}
 \mathbf{s}_4 &= R(\mathbf{s}_3, \mathbf{x}_4) \\
 &= R(R(\mathbf{s}_2, \mathbf{x}_3), \mathbf{x}_4) \\
 &= R(R(R(\mathbf{s}_1, \mathbf{x}_2), \mathbf{x}_3), \mathbf{x}_4) \\
 &= R(R(R(R(\mathbf{s}_0, \mathbf{x}_1), \mathbf{x}_2), \mathbf{x}_3), \mathbf{x}_4)
 \end{aligned}$$

2 Comme on peut l'observer ci-dessus, la valeur de \mathbf{s}_4 est calculée à partir
 3 des entrées $\mathbf{x}_{1..4}$.

4 \mathbf{s}_n et $\hat{\mathbf{y}}_n$ peuvent être vus comme des manières de représenter (on dit aussi
 5 d'**encoder**) toute la séquence d'entrée $\mathbf{x}_{1..n}$.

6 C'est lors de l'étape d'apprentissage que les paramètres de R et de O sont
 7 déterminés, de manière que les informations pertinentes pour une prédiction
 8 au temps t soient véhiculées par le vecteur d'état.

9 L'apprentissage des paramètres est réalisé sur le réseau sous forme déplié.

10 Pour cela, il faut définir une fonction de perte et utiliser l'algorithme de
 11 descente du gradient, par exemple, pour calculer les paramètres du réseau qui
 12 permettent de minimiser la fonction de perte sur des données d'apprentissage.

13 C'est par le biais de la minimisation de la fonction de perte que le réseau
 14 apprend ce dont il faut se souvenir du passé, en d'autres termes, ce qui doit
 15 être représenté dans le vecteur d'état.

16 1.1 Accepteurs et Transducteurs

17 1.1.1 Accepteurs

18 Il est possible de faire en sorte que le réseau ne fasse qu'une prédiction,
 19 à l'issue de la lecture de toute la séquence.

20 Lors de l'apprentissage, ce n'est qu'à l'issue de la prédiction finale que
 21 l'on calcule la fonction de perte et que la mise à jour des paramètres est
 22 réalisée.

23 Exemples :

- 24 1. Prédiction de la langue. Contrairement au CNN, on a maintenant des
 25 séquences de longueurs quelconques et c'est à l'issue de la lecture du
 26 dernier caractère qu'une prédiction est réalisée.

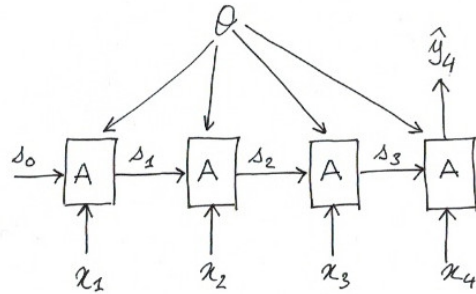


FIGURE 1.3 – un accepteur

2. Prédiction du jugement véhiculé par une phrase. On veut concevoir un réseau qui lit une phrase et prédit, à la fin de la lecture si la phrase véhicule un jugement positif, négatif ou neutre, comme dans l'exemple ci-dessous :

- *J'ai vraiment détesté ce cours* → NEGATIF
- *rarement un cours m'a autant intéressé* → POSITIF
- *Le cours a lieu toutes les semaines en salle 202* → NEUTRE

Le vecteur d'état s_n produit à l'issue du traitement d'une séquence $\mathbf{x}_{1\dots n}$ est souvent utilisée comme entrée d'une autre tâche de prédiction.

On peut imaginer par exemple que s_n constitue l'entrée d'un autre réseau R qui peut être, par exemple, MLP ou un CNN, comme l'illustre la figure 1.4.

s_n correspond à une **représentation** de la séquence $\mathbf{x}_{1\dots n}$, on dit aussi son **encodage**.

Le RNN est appelé dans ce cas là un **encodeur**.

La prédiction est réalisée par R et c'est donc à la sortie de R que la fonction de perte est calculée et que la rétro-propagation est réalisée, pour mettre à jour les paramètres du RNN.

1.1.2 Transducteurs

De manière générale, un transducteur est un programme qui lit une séquence et en produit une autre.

On peut imaginer par exemple un programme qui effectue une traduction d'une langue source vers une langue cible. Il lit une séquence dans la langue

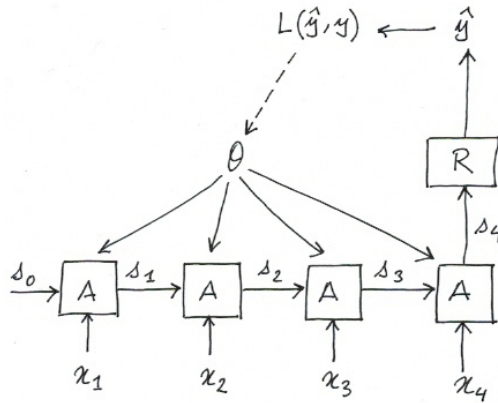


FIGURE 1.4 – un encodeur

1 source et produit une séquence dans la langue cible.

2 Une transduction peut être réalisée par un RNN qui produit une sortie $\hat{\mathbf{y}}_t$
 3 pour toute entrée \mathbf{x}_t .

4 Ceci permet de calculer une perte pour chaque sortie $\hat{\mathbf{y}}_t$ en la comparant
 5 avec la sortie correcte \mathbf{y}_t : $L(\hat{\mathbf{y}}_t, \mathbf{y}_t)$

6 La perte pour toute la séquence peut alors être calculée à partir des pertes
 7 locales, en faisant, par exemple, la somme :

$$L(\hat{\mathbf{y}}_{1\dots n}, \mathbf{y}_{1\dots n}) = \sum_{t=1}^n L(\hat{\mathbf{y}}_t, \mathbf{y}_t)$$

8 Exemple :

9 On veut attribuer une partie de discours (catégorie grammaticale), tel
 10 que : *Verbe*, *Nom*, *Adj* aux différents mots d'un texte, comme dans l'exemple
 11 suivant :

12

\mathbf{x}	Marie	donne	le	la
$\hat{\mathbf{y}}$	Nom	Verbe	Det	Nom

13

14 La prédiction de la partie de discours d'un mot dépend bien sûr du mot
 15 (*manger* est un verbe), mais aussi du contexte du mot (le mot *la* dans notre
 16 exemple est un nom car il est précédé d'un déterminant).

17 C'est le contexte du mot qui est représenté par le vecteur d'état.

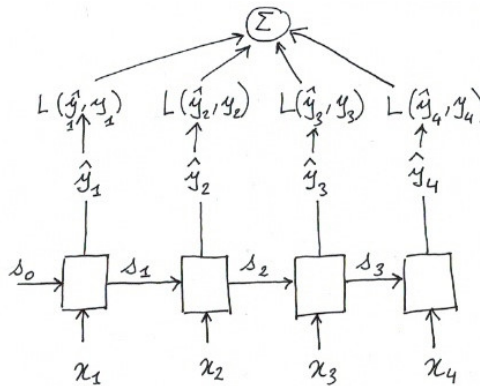


FIGURE 1.5 – un transducteur

1.2 Quelques types de réseaux récurrents

Le modèle de RNN présenté jusque là est très général, il est défini par les deux équations suivantes :

$$\begin{aligned} \mathbf{s}_t &= R(\mathbf{x}_t, \mathbf{s}_{t-1}) \\ \hat{\mathbf{y}}_t &= O(\mathbf{s}_t) \end{aligned}$$

Selon la forme que prennent R et O , différents types de RNN peuvent être définis.

On en verra deux :

- Les RNN simples, notés SRNN, ou réseaux de Elman
- Les *Gated Recurrent Units*, notés GRU

1.2.1 RNN simples

Il s'agit d'un des modèles récurrents les plus simples.

Il est défini de la façon suivante :

$$\begin{aligned} \mathbf{s}_t &= R_{SRNN}(\mathbf{x}_t, \mathbf{s}_{t-1}) \\ &= g(W^s \mathbf{s}_{t-1} + W^x \mathbf{x}_t + \mathbf{b}) \\ \hat{\mathbf{y}}_t &= O_{SRNN}(\mathbf{s}_t) \\ &= \mathbf{s}_t \end{aligned}$$

1 où g est une fonction d'activation non linéaire.
 2 Comme on peut le voir :
 3 — l'état courant \mathbf{s}_t est calculé à partir de l'état précédent \mathbf{s}_{t-1} et de
 4 l'entrée courante \mathbf{x}_t en multipliant \mathbf{s}_{t-1} par la matrice W^s et \mathbf{x}_t par
 5 la matrice W^x puis en faisant la somme des deux vecteurs produits et
 6 d'un vecteur de biais \mathbf{b} ;
 7 — La sortie au temps t n'est autre que l'état \mathbf{s}_t .
 8 Le modèle SRNN est confronté au problème de la disparition du gradient
 9 (*vanishing gradient* en anglais) : lors de la rétro-propagation de l'erreur, cette
 10 dernière diminue rapidement (car elle est multipliée successivement par des
 11 poids dont la valeur est souvent comprise entre 0 et 1).
 12 Cette caractéristique empêche la bonne prise en compte de dépendances
 13 entre observations éloignées dans le temps (le réseau oublie trop vite).
 14 Pour pallier ce problème, on introduit un système de **porte** ou de **filtre**
 15 (*gate* en anglais).
 16 Un filtre \mathbf{f} se présente comme un vecteur composé de 0 et de 1.
 17 Il est combiné à un autre vecteur \mathbf{x} à l'aide d'une opération appelée
 18 **produit de Hadamard** : $\mathbf{f} \circ \mathbf{x}$.
 19 Cette opération consiste simplement à effectuer le produit terme à terme
 20 de deux vecteurs (la composante i de \mathbf{f} est multipliée par la composante i de
 21 \mathbf{x}), comme dans l'exemple suivant :

$$\begin{bmatrix} 4 \\ 5 \\ 7 \\ 9 \end{bmatrix} \circ \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 4 \\ 0 \\ 0 \\ 9 \end{bmatrix}$$

22 Dans cet exemple, le filtre a permis "d'oublier" les composantes 2 et 3 du
 23 vecteur \mathbf{x} .
 24 Un filtre permet donc de décider ce que l'on garde et ce que l'on oublie
 25 dans un état \mathbf{s}_i .
 26 Les filtres offrent un moyen simple de combiner au sein d'un vecteur d'état
 27 \mathbf{s}_i le contenu de l'entrée \mathbf{x}_i et du vecteur d'état \mathbf{s}_{i-1} de la manière suivante :

$$\mathbf{s}_i = \mathbf{f} \circ \mathbf{x}_i + (1 - \mathbf{f}) \circ \mathbf{s}_{i-1}$$

28 Ainsi, on garde les composantes de l'entrée courante \mathbf{x}_i qui correspondent
 29 à des 1 dans le filtre et on oublie les autres et on garde les composantes de
 30 l'état précédent \mathbf{s}_{i-1} qui correspondent à des 0 dans le filtre.

1 Exemple :

$$\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \circ \begin{bmatrix} 10 \\ 11 \\ 12 \\ 13 \\ 14 \\ 15 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 1 \\ 0 \end{bmatrix} \circ \begin{bmatrix} 8 \\ 9 \\ 3 \\ 7 \\ 5 \\ 8 \end{bmatrix} = \begin{bmatrix} 8 \\ 11 \\ 3 \\ 7 \\ 5 \\ 15 \end{bmatrix}$$

2 L'idée, bien entendu, est d'apprendre les filtres lors de l'apprentissage :
3 le réseau apprend l'information qu'il doit retenir de l'état précédent.

4 Cela pose un problème technique : l'algorithme de rétro-propagation im-
5 pose que toutes les fonctions qui participent à la fonction de perte soient
6 dérivables, ce qui n'est pas le cas des filtres dont les composantes ne peuvent
7 prendre qu'une des deux valeurs 0 ou 1.

8 Une solution consiste à approximer le comportement des filtres binaires à
9 l'aide d'une fonction dérivable pour donner naissance à des **filtres dérivables**.

10 Pour cela, on autorise les composantes du filtre à prendre pour valeur
11 n'importe quel réel, ce qui constitue un filtre \mathbf{f}' à valeurs réelles.

12 Ce filtre est ensuite transformé à l'aide de la fonction sigmoïde (σ), ce
13 qui va avoir pour effet de ramener toutes les valeurs à l'intervalle $[0, 1]$ et
14 "repousser" la majorité des valeurs qui composent le vecteur vers 0 ou 1,
15 comme on peut le voir dans la figure 1.6.

16 Le filtre $\sigma(\mathbf{f}')$ est par conséquent composé de valeurs proches de 0 et de
17 1.

18 Lorsqu'on réalise le produit $\sigma(\mathbf{f}') \circ \mathbf{x}$, les composantes de \mathbf{x} correspondant
19 à des valeurs proches de 1 dans $\sigma(\mathbf{f}')$ sont conservées alors que les autres sont
20 oubliées.

21 Ce mécanisme de filtre dérivable est utilisé dans deux types de RNN, les
22 GRU et les LSTM.

23 A chaque étape, de tels réseaux décident, grâce à des filtres dérivables,
24 quelles informations de l'entrée vont être prises en compte dans l'état courant
25 et quelles informations de l'état précédent vont être oubliées.

26 1.2.2 Les Gated Recurrent Units

27 Les GRU mettent en œuvre deux filtres :

28 — Le filtre *update gate*, noté \mathbf{z}_t , qui contrôle quelles informations conte-
29 nues dans \mathbf{s}_{t-1} vont être transmises à \mathbf{s}_t .

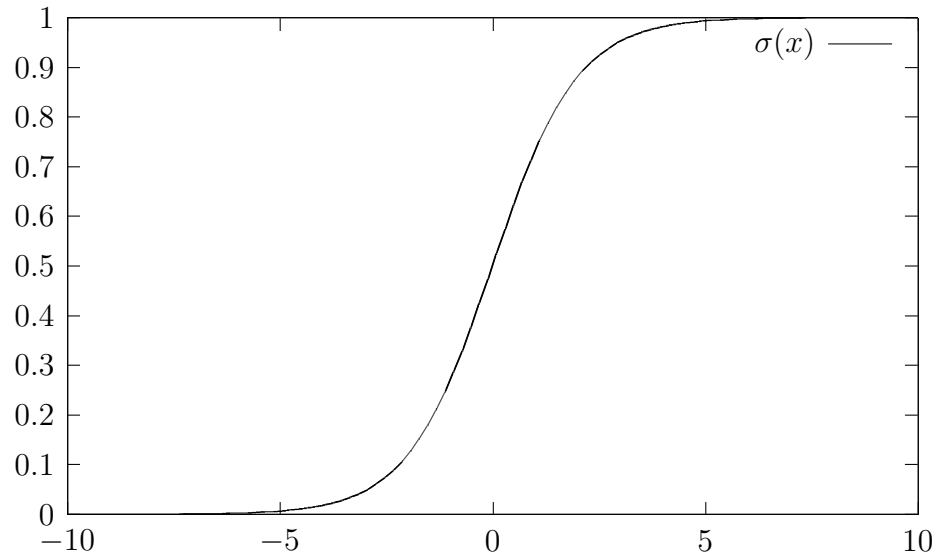


FIGURE 1.6 – La fonction sigmoïde

- 1 — Le filtre *reset gate*, noté \mathbf{r}_t , qui contrôle quelles informations contenues
 2 dans \mathbf{s}_{t-1} vont être combinées à \mathbf{x}_t .

3 **Calcul de \mathbf{z}_t**

4 Le filtre \mathbf{z}_t est calculé à partir de :

- 5 — l'entrée courante : \mathbf{x}_t
 6 — l'état précédent : \mathbf{s}_{t-1}
 7 — deux matrices de paramètres X^z et U^z

$$\mathbf{z}_t = \sigma(W^z \mathbf{x}_t + U^z \mathbf{s}_{t-1})$$

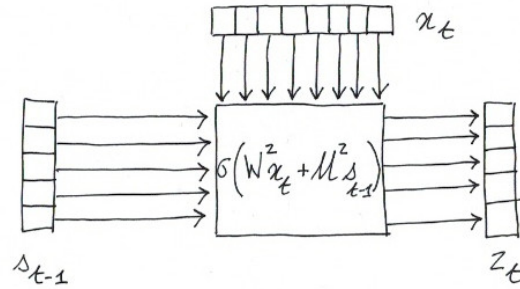
8 Ce calcul est représenté schématiquement dans la figure 1.7.

9 Lors de la prédiction, les matrices W^z et U^z sont constantes (elles ont été
 10 calculées lors de la phase d'apprentissage) mais pour tout couple $(\mathbf{x}_t, \mathbf{s}_{t-1})$,
 11 une nouvelle valeur du filtre est calculée.

12 **Calcul de \mathbf{r}_t**

13 Le calcul de \mathbf{r}_t suit le même principe que celui de \mathbf{z}_t , avec les matrices
 14 X^r et U^r :

$$\mathbf{r}_t = \sigma(W^r \mathbf{x}_t + U^r \mathbf{s}_{t-1})$$

FIGURE 1.7 – Représentation schématique du calcul du filtre z_t 1 **Calcul de s_t** 2 L'équation générale de calcul de s_t est :

$$s_t = R(s_{t-1}, x_t)$$

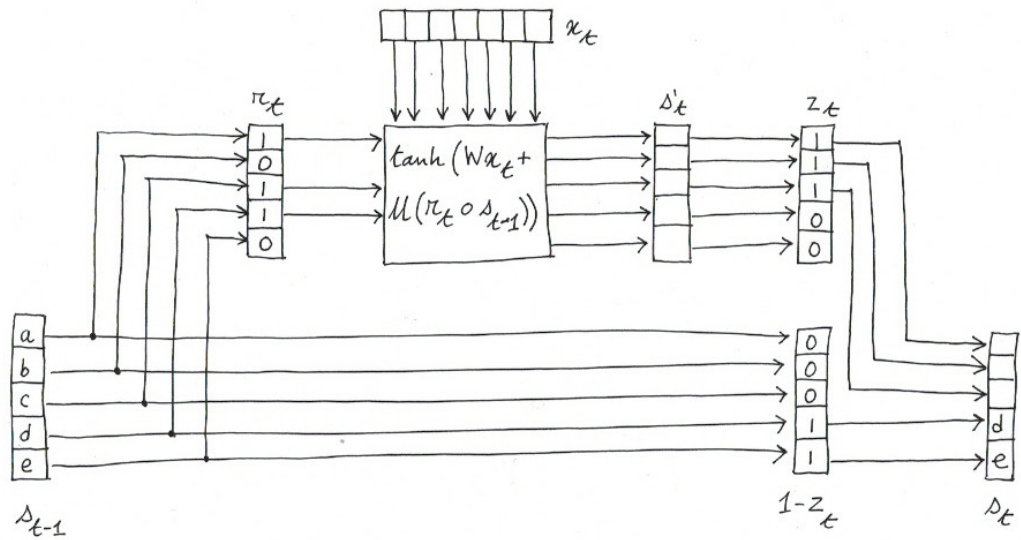
3 Dans le cas des GRU, cette équation prend la forme suivante :

$$s_t = z_t \circ s_{t-1} + (1 - z_t) \circ s'_t$$

4 s'_t est appelé l'état interne au temps t , il est calculé de la manière sui-
5 vante :

$$s'_t = \tanh(Wx_t + U(r_t \circ s_{t-1}))$$

6 où \tanh est la fonction tangente hyperbolique (voir chapitre sur les fonc-
7 tions).8 Ce calcul est représenté schématiquement dans la figure 1.8 qui donne
9 une idée générale du comportement du GRU.10 Afin de ne pas surcharger le schéma, nous avons omis le calcul des filtres
11 z_t et r_t à partir de l'entrée x_t et du vecteur d'état s_{t-1} .12 On observe, dans la partie basse du schéma, que les composantes 4 et 5,
13 correspondant à d et e , du vecteur d'état s_{t-1} sont recopiés dans le vecteur
14 d'état s_t .15 Cette recopie est réalisée par la combinaison du vecteur d'état s_{t-1} et du
16 filtre $1 - z_t$.17 Le vecteur d'état s_{t-1} est aussi combiné, dans la partie haute du schéma,
18 avec le filtre r_t , qui ne garde du premier que les composantes 1, 3 et 4 (les
19 valeurs a , c et d) et oublie par conséquent les valeurs b et e .

FIGURE 1.8 – Représentation schématique du calcul du vecteur d'état \mathbf{s}_t

- 1 Les composantes 1, 3 et 4 sont combinées avec l'entrée \mathbf{x}_t , à l'aide des
- 2 matrices W et U pour former l'état interne \mathbf{s}'_t . Ce dernier est alors combiné
- 3 avec le filtre \mathbf{z}_t , qui n'en garde que les trois premières composantes.
- 4 Le vecteur d'état \mathbf{s}_t est alors construit, en gardant les trois premières
- 5 composantes de \mathbf{s}'_t et les deux dernières composantes de \mathbf{s}_{t-1} .