

Amlioration de la gnration de lgendes d'images avec RBM sur le module de captioning de BLIP (blip-image-captioning-base)

Fortun HOUESSOU

Sous la supervision de M. Jrme Lacaille

Mars 2025

Abstract

Le Image captioning, qui consiste gnrer des descriptions textuelles partir de contenus visuels, a connu des avances majeures grce aux modles Vision-Language comme BLIP. Cependant, ces approches rencontrent encore des difficults produire des lgendes la fois prcises et contextuellement cohrentes, notamment pour des images complexes ou ambigus. Dans ce travail, nous proposons une amlioration innovante en intgrant une Machine de Boltzmann Restreinte (RBM) dans le processus de gnration de lgendes de BLIP. Les RBM, reconnues pour leur capacit apprendre des representations latentes profondes, permettent d'affiner la cohrence smantique des lgendes gnres sans ncessiter un reentrainement massif sur de grandes bases de donnees. Nous explorons diffrentes stratgies d'intgration, o la RBM intervient pour rvaluer ou ajuster les sorties de BLIP, afin d'amliorer la pertinence et la richesse des descriptions. Des experiences approfondies menes sur des benchmarks standard montrent que cette approche hybride permet de gnrer des lgendes plus naturelles et contextuellement adaptes, ouvrant ainsi de nouvelles perspectives pour l'amlioration des modles Vision-Language.

Mots-cl : Captioning d'images, Modles Vision-Langage (VLM), BLIP, Machine de Boltzmann Restreinte (RBM), Apprentissage profond.

Contents

1	Introduction	3
2	Image - Captioning : tat de l'Art	3
2.1	Approches classiques du captioning d'images	3
2.1.1	Rseaux de Neurones Rcurrents (RNN)	3
2.1.2	Transformers	4
2.1.3	Modles Vision-Langage	4
3	Vision-Language Models (VLM)	4
3.1	Prtraitement et Encodage	6
3.2	Fusion des Représentations : Alignement Visuel-Linguistique	6

4	BLIP (Bootstrapping Language-Image Pre-training)	9
4.0.1	Fonctionnement gnral	10
4.0.2	Architecture globale	11
5	Blip-Image-Captioning-Base	11
6	La Machine de Boltzmann Restreinte (RBM)	12
6.1	Pourquoi la machine de Boltzmann restreinte ?	13
6.2	Structure d'une RBM	13
6.3	Fonction nergetique	14
6.4	La RBM comme modle probabiliste	14
6.5	Formalisme	15
7	Quest-ce qu'une RBM Gaussienne-Binaire?	18
8	Fonctionnement d'une RBM-GB	18
8.1	Activation des neurones caches (identique la RBM Binaire-Binaire) . . .	19
8.2	Activation des neurones visibles (modifi pour les valeurs continues) . . .	19
9	Entranement d'une RBM gaussienne	19
10	Une fois entrane, comment on utilise la RBM gaussienne ?	20
11	Limites d'une RBM gaussienne	20

1 Introduction

L'association entre la vision par ordinateur et le traitement automatique du langage naturel a révolutionné la manière dont les machines interprètent et décrivent le contenu visuel. Cette convergence a donné naissance à des modèles Vision-Language Models (VLMs), capables d'analyser une image et de générer une description textuelle cohérente et pertinente. Parmi les modèles les plus performants dans ce domaine est le Bootstrapping Language-Image Pre-training (BLIP), qui exploite des techniques avancées d'apprentissage auto-supervisé pour améliorer la compréhension des relations entre les images et le texte.

Cependant, malgré son efficacité, BLIP souffre de certaines limitations, notamment une tendance à produire des descriptions génériques, parfois trop peu informatives ou manquant de diversité sémantique. Cette limitation est due en partie aux biais inhérents aux données d'entraînement et aux mécanismes d'optimisation utilisés. Pour remédier à ces faiblesses, il est crucial d'explorer des méthodes complémentaires qui puissent affiner les représentations latentes du modèle et améliorer la richesse du langage généré.

Dans cette optique, l'intégration d'un modèle probabiliste comme la Machine de Boltzmann Restreinte (RBM) représente une piste prometteuse. Les RBMs sont largement utilisés dans l'apprentissage non supervisé et la modélisation des distributions complexes. Elles offrent un moyen d'améliorer les représentations cachées en capturant des interactions non linéaires subtiles dans les données. En les combinant avec le module de captioning de BLIP, nous visons à enrichir la qualité des descriptions générées en affinant les représentations sémantiques sous-jacentes.

Ce travail explore donc l'intégration d'une RBM au sein du pipeline de génération de descriptions d'images de BLIP, en évaluant son impact sur la précision, la diversité et la pertinence linguistique des légendes produites. Cette approche pourrait ouvrir de nouvelles perspectives pour l'amélioration des modèles Vision-Language et leur application à des domaines exigeant une compréhension fine du contenu visuel, comme l'accessibilité numérique ou l'annotation automatique d'images.

2 Image - Captioning : l'état de l'Art

Le **captioning d'images**, ou annotation automatique d'images, est une tâche en vision par ordinateur et de langage consistant à générer une description textuelle cohérente et précise du contenu d'une image. Cette tâche est complexe car elle nécessite une compréhension approfondie des objets présents, de leurs actions, de leurs interactions et du contexte global de l'image.

2.1 Approches classiques du captioning d'images

Plusieurs approches ont été explorées parmi lesquelles :

2.1.1 Réseaux de Neurones Récurrents (RNN)

Les premières approches utilisaient des RNN, notamment des LSTM (Long Short-Term Memory), pour générer des descriptions séquentielles des images. Un pipeline typique impliquait l'utilisation d'un réseau de neurones convolutifs (CNN) pour extraire des caractéristiques visuelles de l'image, suivie d'un RNN pour générer la légende basée sur ces caractéristiques.

tiques. Par exemple, l'article *Language Models for Image Captioning: The Quirks and What Works* explore l'utilisation des RNN pour cette tche. [?].

2.1.2 Transformers

Les modles *Transformers*, introduits initialement pour le traitement du langage naturel, ont t adapts au captioning d'images en raison de leur capacit modliser des dpendances longue porte dans les donnes squentielles. Ces modles utilisent des mcanismes d'attention pour se concentrer sur diffrentes parties de l'image lors de la gnration de chaque mot de la lgende. L'article *A Review of Transformer-Based Approaches for Image Captioning* fournit une revue exhaustive des modles de captioning d'images bass sur les Transformers[?].

2.1.3 Modles Vision-Langage

Les modles *Vision-Langage* intgrent des informations visuelles et textuelles pour amliorer la performance du captioning d'images. Ces modles sont pr-entrans sur de grandes quantits de donnes d'images et de textes associs, puis affins pour des tches spcifiques comme le captioning. Le modle blip-image-captioning [8] qui fera l'objet de notre tude est un exemple concret des performances atteintes avec les VLMs.

Mais c'est quoi rellement un VLM ?

3 Vision-Language Models (VLM)

Les modles de Vision Langage (VLM) sont des architectures multimodales qui combinent le traitement visuel et linguistique pour accomplir des tches varies telles que la gnration de lgendes d'images, la recherche d'images par texte, ou l'association image-texte. Ces modles ont gagn en popularit grce leur capacit fusionner les informations textuelles et visuelles, permettant ainsi de rsoudre des problmes complexes en vision par ordinateur et en traitement du langage naturel. Parmi les modles les plus connus dans ce domaine figurent CLIP [11] et BLIP [8].

BLIP, en particulier, adopte une approche de pr-entrainement utilisant des donnes web massives et combine un encodeur d'image avec un dcodeur de texte pour gnrer des descriptions adaptes aux images. Cette approche a dmontr son efficacit en gnration de lgendes et interprétation multimodale. De mme, CLIP [11] se distingue par sa capacit associer des images et des descriptions textuelles, facilitant ainsi des applications comme la recherche d'images par texte et l'analyse d'images partir de requetes textuelles.

Les VLMs sont galement prsents dans des travaux comme VisualBERT [9], qui fusionne les informations visuelles et textuelles dans un modle transformeur pour des tches de comprhension multimodale, telles que la rponse des questions visuelles. Ce modle a montr qu'il est possible d'aligner les reprsentations visuelles et textuelles dans un espace partag pour des performances solides sur des tches de vision et de langage.

Un autre modle significatif dans ce domaine est ViLT [7], qui a introduit une approche simplifie, ne recourant ni aux convolutions ni une supervision rgionale, permettant de traiter simultanment les images et les textes tout en rduisant la complexit du modle. Cela a ouvert la voie des applications plus lgres tout en maintenant des performances comptitives.

Dans le cadre de la conduite autonome et d'autres applications de sécurité, des travaux comme VLM-RL [6] ont exploré l'intégration des VLMs avec l'apprentissage par renforcement pour améliorer la prise de décision basée sur des informations multimodales, notamment en gestion de scénarios visuels et textuels complexes.

Les capacités de raisonnement spatial sont également un domaine d'amélioration important pour les VLMs, comme le montre SpatialVLM [3], qui cherche à doter ces modèles de la capacité d'effectuer des raisonnements géométriques et spatiaux, ce qui est essentiel pour des applications telles que la navigation autonome et la manipulation d'objets en 3D.

Enfin, des modèles comme VILA [5] et Xmodel-VLM [10] se concentrent sur des techniques de pré-entraînement et d'optimisation pour rendre les VLMs plus efficaces et adaptés aux environnements de production, en réduisant la taille des modèles tout en préservant leur capacité à traiter des informations visuelles et textuelles de manière précise.

Dans la suite de ce rapport sur l'état de l'Art des VLMs, de Blip et RBM, nous commençons par présenter de façon générale, le fonctionnement d'un VLM puis le rôle des mathématiques derrière ces modèles, et enfin nous traiterons les différents modèles qui feront l'objet de ce sujet savoir BLIP (Bootstrapping Language-Image Pre-training) et RBM (Restricted Boltzmann machine).

Les **modèles de Vision et Langage** (VLMs) sont des systèmes qui essaient de comprendre à la fois des **images** et des **mots**. Imaginons qu'on montre une image de chat à un enfant et qu'on lui dise "C'est un chat". Le cerveau de l'enfant va faire le lien entre l'image du chat et le mot "chat". C'est exactement ce que font les VLMs, mais d'une manière beaucoup plus complexe et dans un ordinateur !

Comment ça marche concrètement ?

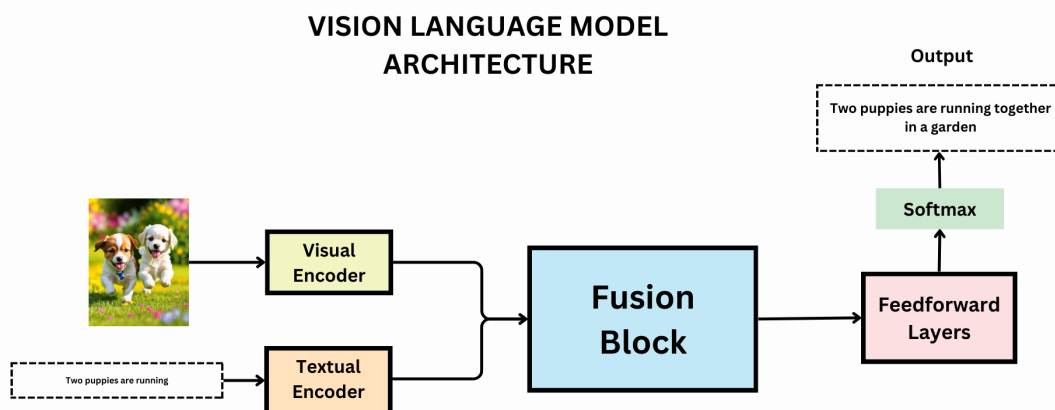


Figure 1: VLM

Le **VLM** combine la compréhension des **images** et des **mots**. Il analyse d'abord l'image pour en extraire des informations visuelles (formes, objets, etc.) via un réseau de neurones. Ensuite, les mots sont transformés en représentations que le modèle peut comprendre.

Enfin, le VLM fusionne les informations des images et des mots pour accomplir des tâches comme l'association image-texte ou la génération de descriptions d'images. L'objectif est d'apprendre à lier efficacement ces deux modalités pour mieux comprendre et interagir avec le monde visuel et linguistique simultanément. C'est comme un Homme qui entend le mot "chat" et sait un peu quelle forme s'attendre, parce que des chats il en a beaucoup vu, de différentes couleurs et dans différents horizons.

Un peu de formalisme :

L'entraînement du modèle BLIP s'effectue en plusieurs étapes :

3.1 Prtraitement et Encodage

Soit x une entrée visuelle (image) et t une entrée textuelle (texte). Dans la plupart des VLMs, l'image x est d'abord passée à travers un encodeur d'image f_{image} , typiquement un réseau de neurones convolutifs (CNN) ou un modèle transformeur ou vision encodeur adapté à l'image (par exemple, Vision Transformer, ViT). L'objectif est de projeter l'image dans un espace de représentation latente $\mathcal{Z}_{\text{image}} \in R^d$.

$$\mathcal{Z}_{\text{image}} = f_{\text{image}}(x)$$

De même, le texte t est encodé par un modèle de langage, tel que le Transformer, en une représentation vectorielle $\mathcal{Z}_{\text{texte}} \in R^d$. Le processus de traitement du texte suit un encodage similaire avec une fonction f_{texte} , qu'on ne connaît pas a priori (BERT par exemple).

$$\mathcal{Z}_{\text{texte}} = f_{\text{texte}}(t)$$

Dans le processus d'embedding, certains modèles de VLM utilisent du : self-attention.

3.2 Fusion des Représentations : Alignement Visuel-Linguistique

Pour regrouper ces deux embeddings dans un même espace latent, il y a deux grandes stratégies :

A. Projection directe

On passe $\mathcal{Z}_{\text{image}}$ et $\mathcal{Z}_{\text{texte}}$ dans des Multi-Layer Perceptrons (MLP) [4] ou des couches linéaires afin de les projeter dans un **même espace latent**. Une fois cette projection effectuée, on calcule simplement la distance entre ces deux représentations, par exemple via une similarité :

$$\text{similarité}(\mathcal{Z}_{\text{image}}, \mathcal{Z}_{\text{texte}})$$

Les fonctions de similarité font l'objet d'une estimation de la distance entre deux embeddings. C'est une façon de représenter par un nombre les différences ou rapprochements entre entités.

Fonctions de similarité

Bases sur la distance ou l'angle

- **Similarité cosinus (Cosine Similarity) :**

$$\text{sim}(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$$

- **Distance euclidienne (Euclidean Distance) :**

$$d(x, y) = \sqrt{\sum_i (x_i - y_i)^2}$$

Souvent transformé en :

$$\text{sim}(x, y) = \frac{1}{1 + d(x, y)}$$

Bases sur des probabilités

- **Divergence de Kullback-Leibler (KL Divergence) :**

Pour deux distributions de probabilités P et Q :

$$D_{KL}(P \parallel Q) = \sum_i P(i) \log \left(\frac{P(i)}{Q(i)} \right)$$

Pas symétrique et non bornée.

- **Divergence de Jensen-Shannon (JS Divergence) :**

Variante symétrique et plus stable du KL :

$$D_{JS}(P, Q) = \frac{1}{2} D_{KL}(P \parallel M) + \frac{1}{2} D_{KL}(Q \parallel M)$$

où

$$M = \frac{1}{2}(P + Q)$$

Bases sur des ensembles

- **Similarité de Jaccard (Jaccard Similarity) :**

Pour des ensembles A et B :

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

- **Coefficient de recouvrement (Overlap Coefficient) :**

$$\text{Overlap}(A, B) = \frac{|A \cap B|}{\min(|A|, |B|)}$$

C'est précisément l'approche adoptée par le modèle CLIP.

B. Cross-Attention

Dans certains modèles comme BLIP, avant la projection des données bimodales (image/texte) dans l'espace latent, un mécanisme d'attention croisée (cross-attention) est utilisé dans les architectures de type transformer. Ce mécanisme permet de calculer une attention mutuelle entre les représentations des images et du texte afin d'extraire des correspondances pertinentes entre les modalités. Les articles [2] et [1] introduisent et expliquent dans les détails le mécanisme d'attention. Le calcul de l'attention croisée peut être formalisé comme suit :

$$\text{Attn}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

o Q (Query) est la requête (environnement textuel ou visuel), K (Key) est la clé (respectivement l'information visuelle ou textuelle), et V (Value) est la valeur : [14] et [13]. Ce mécanisme est appliqué pour apprendre les relations entre $\mathcal{Z}_{\text{image}}$ et $\mathcal{Z}_{\text{texte}}$, en ajustant les poids des connexions neuronales selon l'importance de chaque élément textuel ou visuel pour le modèle.

Optimisation et Apprentissage

L'entraînement des VLMs repose sur la minimisation d'une fonction de perte \mathcal{L} qui permet au modèle d'apprendre à associer de manière optimale les informations visuelles et textuelles. Par exemple, une fonction de perte typique dans le cas d'une tâche de classification d'images par texte est la suivante :

$$\mathcal{L} = \sum_i [\log P(y_i | \mathcal{Z}_{\text{image}}) + \log P(y_i | \mathcal{Z}_{\text{texte}})]$$

o y_i est l'étiquette associée à l'entrée x_i ou t_i . Cette formulation implique une double minimisation des erreurs sur les prédictions visuelles et textuelles.

Un autre type de perte, utilisé dans les modèles de type *contrastive learning* (comme CLIP [11]), repose sur l'idée de maximiser la similarité entre les représentations textuelles et visuelles d'une même instance tout en minimisant la similarité entre les paires d'images et de textes différentes. La fonction de perte contraste peut être formulée ainsi :

$$\mathcal{L}_{\text{contrastive}} = - \sum_i \log \frac{\exp(\text{sim}(\mathcal{Z}_{\text{image}}^i, \mathcal{Z}_{\text{texte}}^i))}{\sum_j \exp(\text{sim}(\mathcal{Z}_{\text{image}}^i, \mathcal{Z}_{\text{texte}}^j))}$$

o $\text{sim}(a, b)$ représente la similarité entre les vecteurs a et b , souvent mesurée par un produit scalaire ou un cosinus.

Architecture générale d'un VLM (Vision-Language Model)

1 Embedding

- Image \rightarrow Patches (ViT) ou CNN \rightarrow Embeddings visuels.
- Texte \rightarrow Tokenisation \rightarrow Embeddings textuels.

2 Self-Attention (selon le modèle)

- Appliqués sur les embeddings :
- Image \rightarrow Self-attention entre patches.
- Texte \rightarrow Self-attention entre tokens.

3 Cross-Attention (si fusion, cas de BLIP)

- Texte attend sur image ou inversement.
- Q = texte, K, V = image (ou l'inverse selon le modèle).
- Permet d'associer finement les informations. (On dit que le texte et l'image se regardent)

4 Projection

- Alignement dans un même espace vectoriel via couches linéaires.
- Image \rightarrow projection.
- Texte \rightarrow projection.

5 Similarité

- Calcul par produit scalaire ou cosin similarity entre vecteurs projetés.

6 Contrastive Loss

- But :
 - Maximiser la similarité des bonnes paires (image, texte liés).
 - Minimiser celle des paires incorrectes.

””

4 BLIP (Bootstrapping Language-Image Pre-training)

BLIP (*Bootstrapping Language-Image Pre-training*) [8] est un framework conçu pour améliorer les performances sur les tâches vision-langage. Il se distingue par sa capacité à exploiter la fois des **données brutes web-scale** et des annotations de qualité, grâce à une architecture flexible et un apprentissage par bootstrapping.

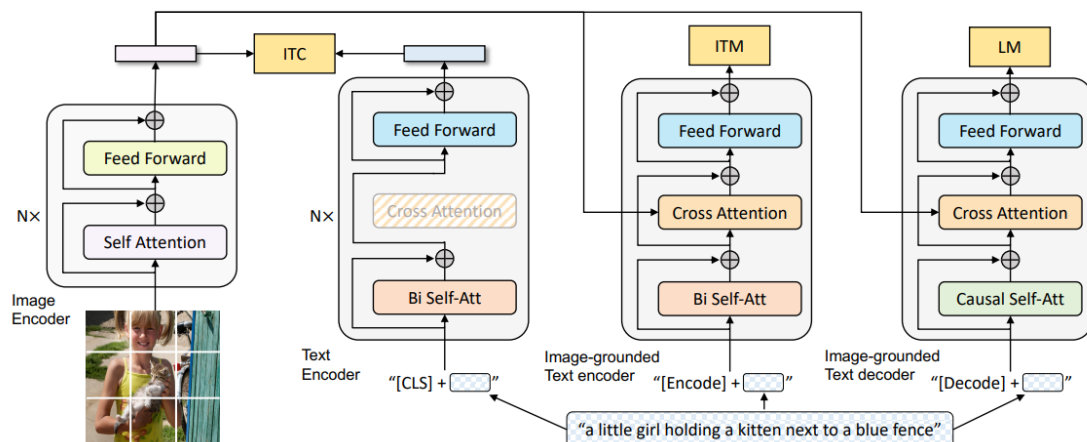


Figure 2: BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation

4.0.1 Fonctionnement gnral

BLIP repose principalement sur trois tches de pr-entrainement complmentaires :

- **Image-Text Contrastive Learning (ITC)** : maximise la similarit entre les representations dune image et de sa lgende correspondante, et minimise celles des paires non correspondantes.

L’ITC est une fonction de perte utilise pour aligner les representations des images et des textes dans un espace commun. Elle est inspire de la InfoNCE loss et fonctionne comme une perte contrastive entre les paires image-texte positives et ngatives.

$$\mathcal{L}_{ITC} = -\frac{1}{N} \sum_{i=1}^N \left(\log \frac{e^{\text{sim}(z_i^I, z_i^T)/\tau}}{\sum_{j=1}^N e^{\text{sim}(z_i^I, z_j^T)/\tau}} \right)$$

o z_i^I et z_i^T sont les embeddings dimage et de texte, sim est la similarit cosinus, τ est la temprature (temperature scaling) et N la taille du lot des donnes traites.

- **Image-Text Matching (ITM)** : discrimine si une paire image-texte est correcte ou incorrecte via une classification binaire, souvent couple une fonction de perte d’entropie croise classique.

$$\mathcal{L}_{ITM} = -\frac{1}{N} \sum_{i=1}^N [y_i \log P(y_i|I_i, T_i) + (1 - y_i) \log(1 - P(y_i|I_i, T_i))]$$

- **Image-Conditioned Language Modeling (LM)** : gnre du texte conditionnellement une image.

$$\mathcal{L}_{LM} = -\frac{1}{N} \sum_{i=1}^N \sum_{t=1}^{T_i} \log P(w_t^i | w_{<t}^i, I_i)$$

o :

- N est le nombre d’exemples dans le batch.
- T_i est la longueur de la squence de texte pour l’exemple i . (nombre de token)
- w_t^i est le mot la position t dans la squence i .
- $w_{<t}^i$ represente tous les mots prcdents dans la squence.
- I_i est l’image associe.
- $P(w_t^i | w_{<t}^i, I_i)$ est la probabilit prdite du mot w_t^i tant donn l’image et les mots prcdents.

4.0.2 Architecture globale

BLIP s'appuie sur une combinaison de :

- **Encoders visuels** (souvent des Vision Transformers - ViT)
- **Encoders textuels** (type BERT)
- **Decoders textuels** (transformer auto-régressif)

Pendant l'entraînement, les modèles apprennent dans un espace latent partagé où les embeddings d'images et de textes sont projetés, permettant une compatibilité sémantique entre les deux modalités.

Dans ce projet, nous nous intéressons plus précisément au modèle BLIP-Captionner du framework BLIP, qui en sortie, génère une description textuelle en réponse à une image reçue en entrée.

5 Blip-Image-Captioning-Base

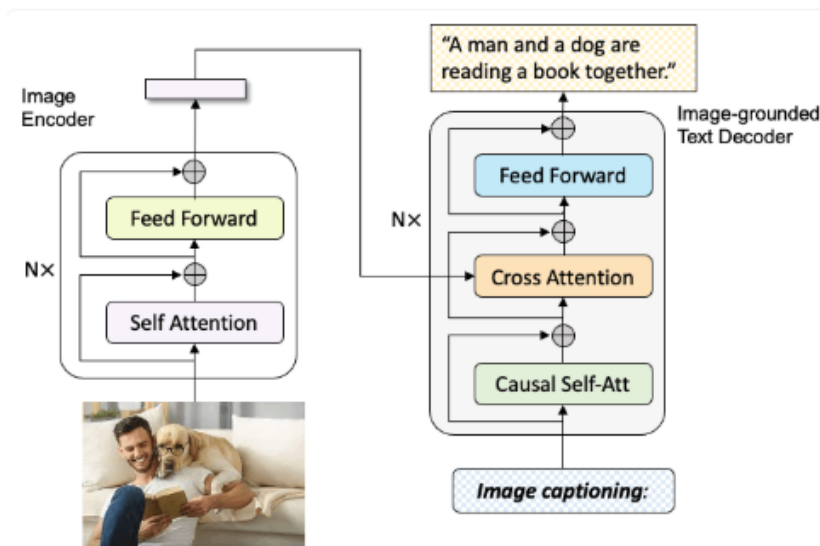


Figure 3: Blip-image-captioning-base

Pour la génération de légende d'image (image captioning) le modèle blip-image-captioning-base utilise :

- Un encodeur d'image (ViT) composé de couches d'auto attention (SA) et d'un réseau Feed Forward (FF) le tout entraîné sur des bases de données annotées par des humains comme Coco et des données brutes du web et nettoyées par bootstrapping via la méthode CapFilt. CapFilt est un ensemble formé d'un Captioner (Image-grounded Text Decoder)[8] et d'un Filter (Image-grounded Text Encoder) qui sont des modules du framework BLIP et qui servent à réannoter (génération de légende) les images brutes obtenues sur le web.

- Le dcodeur de Blip-image-captioning-base est compos d'un rseau FF, pr-entrain avec de l'auto attention bidirectionnelle sur coco dataset. La perte utilise est l'ITC (confre figure 2). Il n'y a pas d'attention croise (le module "Cross Attention" est gel dans un premier temps). Dans un second temps, le meme modle pr-entrain est repris et cette fois ci avec l'attention croise(image-grounded Text encodeur sur la figure 2) et reentrain sur Coco dataset pour une tache de classification binaire [8]. La perte utilise est l'ITM. Enfin le dcodeur de blip-image-captioning est obtenu en remplaant les couches d'auto attention bidirectionnelle par des couches d'auto-attention causale comme sur la figure 2 (image-grounded Text Decoder). La perte utilise lors de l'entrainement est la LM.

Notre objectif est d'augmenter les performances de blip-image-captioning pour affiner la gnration de lgende d'image par insertion d'une RBM (Machine de Boltzmann Restreinte) entre l'encodeur et le dcodeur du modele.

6 La Machine de Boltzmann Restreinte (RBM)

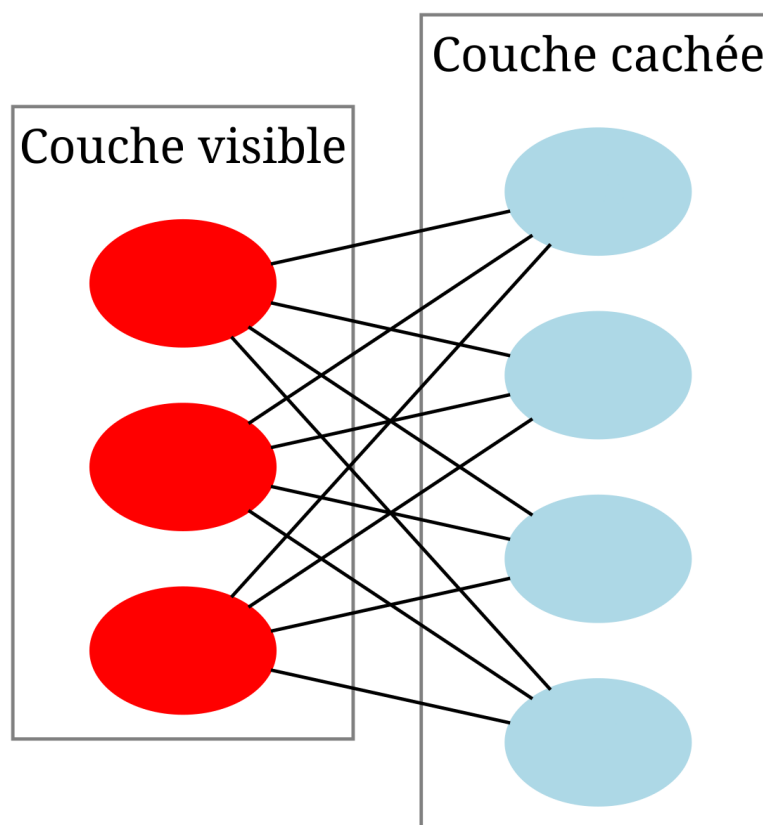


Figure 4: Restricted Boltzmann machine

La Machine de Boltzmann Restreinte (RBM), introduite par [12] et popularise pour l'apprentissage non supervis, est un modle probabiliste gnratif appartenant la famille des modles nergetiques. Apparues dans les annes 1980, les Machines de Boltzmann Restreintes (RBM) ont bnfici d'une meilleure attention avec lessor des capacits de calcul

et leur intégration dans des structures d'apprentissage profondes, en particulier les Deep Belief Networks (DBN). En superposant plusieurs RBM, il devient possible d'extraire des représentations successives des données, passant de caractéristiques élémentaires des concepts plus abstraits. Par ailleurs, lorsqu'un RBM est exploité pour transformer des entrées en représentations latentes, elle peut être assimilée à un réseau de neurones à propagation avant, optimisant ainsi son intégration dans des modèles supervisés pour améliorer l'efficacité de l'apprentissage.

6.1 Pourquoi la machine de Boltzmann restreinte ?

L'utilisation d'une **Restricted Boltzmann Machine (RBM)** pour améliorer l'**image captioning** est un choix intuitif et judicieux, notamment en raison de sa capacité à apprendre des représentations riches et compactes des données visuelles de manière **non supervisée**. En effet, un RBM permet de modéliser efficacement une distribution de probabilité inconnue en découvrant des caractéristiques latentes pertinentes dans les images. Grâce à sa structure bipartite composée de **variables visibles** (les pixels ou des features extraites d'un réseau de neurones) et de **variables latentes** (invisibles), il est capable de capturer des **dépendances complexes** entre les pixels d'une image et d'encoder des relations de haut niveau, facilitant ainsi la génération de descriptions précises et cohérentes.

Dans un cadre d'apprentissage probabiliste, le RBM optimise ses **paramètres** θ via des méthodes comme **l'estimation du maximum de vraisemblance**, ce qui lui permet d'encoder efficacement la structure sous-jacente des images. Cette approche est particulièrement bénéfique pour l'image captioning, où il est crucial de comprendre non seulement les objets présents dans une image, mais aussi leurs relations contextuelles et sémantiques. En apprenant des représentations de manière non supervisée, le RBM permet de généraliser sur de nouvelles images et d'améliorer la robustesse du modèle de génération de descriptions textuelles.

Enfin, en intégrant un RBM en amont d'un modèle de génération de texte, on bénéficie d'une représentation d'entrée plus expressive et informative, ce qui améliore la qualité des légendes générées. Cette synergie entre **modélisation probabiliste** et **apprentissage profond** fait du RBM un choix judicieux et intuitif pour l'amélioration des systèmes d'image captioning.

6.2 Structure d'une RBM

Une RBM est un réseau bipartite composé de :

- **Une couche visible** $\mathbf{v} = (v_1, v_2, \dots, v_m)$ représentant les données observables.
- **Une couche cachée** $\mathbf{h} = (h_1, h_2, \dots, h_n)$ permettant de capturer des dépendances complexes (cachées).

Contrairement aux Machines de Boltzmann standards, il n'y a **pas de connexions intra-couche** :

- Aucune connexion entre les neurones visibles.
- Aucune connexion entre les neurones cachés.

6.3 Fonction nergtique

Le cur de la RBM repose sur une fonction dnergie dfinie par :

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{v}^\top \mathbf{W} \mathbf{h} - \mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top \mathbf{h}$$

o :

- \mathbf{W} est la matrice des poids entre les couches visible et cache.
- \mathbf{b} et \mathbf{c} sont les biais des couches visible et cache respectivement.
- \mathbf{v} et \mathbf{h} sont les tats des neurones sur les couches visible et cache respectivement.

6.4 La RBM comme modle probabiliste

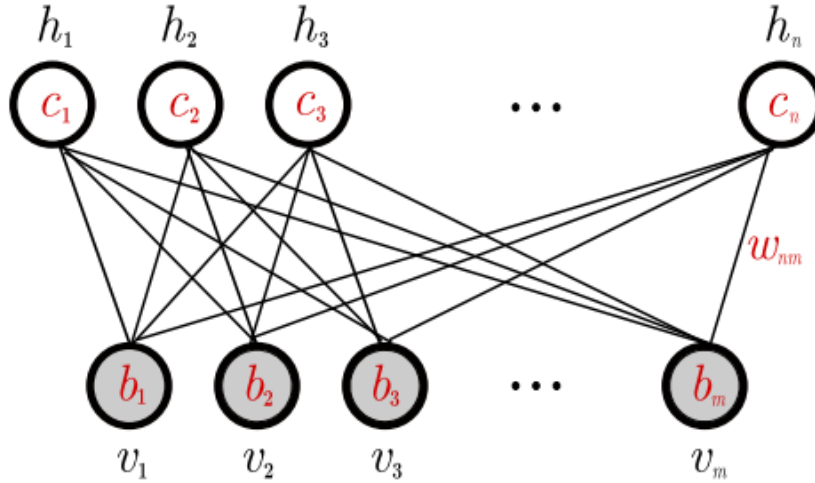


Figure 5: Restricted Boltzmann machine

Une Machine de Boltzmann Restreinte (RBM) est un modle probabiliste bas sur un champ de Markov alatoire structur sous forme dun graphe biparti non orient. Elle est constitue de m units visibles $V = (V_1, \dots, V_m)$, qui representent les donnees observables, et de n units caches $H = (H_1, \dots, H_n)$, charges de capturer les interactions entre ces donnees.

Dans le cadre des RBM binaires, qui sont lobjet de cette tude, les variables alatoires (V, H) prennent des valeurs dans $\{0, 1\}^{m+n}$. La distribution conjointe associe ce modle suit la loi de Gibbs :

$$p(v, h) = \frac{1}{Z} e^{-E(v, h)} \quad (1)$$

avec l'nergie dfinie par :

$$E(v, h) = - \sum_{i=1}^n \sum_{j=1}^m w_{ij} h_i v_j - \sum_{j=1}^m b_j v_j - \sum_{i=1}^n c_i h_i. \quad (2)$$

Ici, pour tout $i \in \{1, \dots, n\}$ et $j \in \{1, \dots, m\}$, w_{ij} dsigne un coefficient rel associ la connexion entre lunit cache H_i et lunit visible V_j . Les paramtres b_j et c_i correspondent respectivement aux biais appliqus aux units visibles et caches.

L'entrainement d'une **Machine de Boltzmann Restreinte (RBM)** repose principalement sur des méthodes d'optimisation adaptées aux modèles probabilistes à base d'énergie. L'algorithme d'apprentissage le plus couramment utilisé est **Contrastive Divergence (CD)**, introduit par Geoffrey Hinton. D'autres algorithmes comme : Parallel Tempering et Persistent Contrastive Divergence (PCD) sont aussi utilisés.

L'apprentissage des paramètres $(\mathbf{W}, \mathbf{b}, \mathbf{c})$ repose sur une estimation approchée du gradient du logarithme de la vraisemblance :

$$\Delta w_{ij} = \epsilon (\langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{model}}) \quad (3)$$

$$\Delta W = \epsilon \cdot (\langle v h^T \rangle_{\text{données}} - \langle v h^T \rangle_{\text{modèle}}) \quad (4)$$

où :

- ϵ est le taux d'apprentissage.
- $\langle \cdot \rangle_{\text{data}}$ est l'espérance sous la distribution empirique des données.
- $\langle \cdot \rangle_{\text{model}}$ est l'espérance sous la distribution modélisée.

Le processus implique des étapes de **Gibbs Sampling**, alternant entre la mise à jour de \mathbf{h} et \mathbf{v} :

$$P(h_j = 1 | \mathbf{v}) = \sigma \left(\sum_i w_{ij} v_i + c_j \right) \quad (5)$$

et

$$P(v_i = 1 | \mathbf{h}) = \sigma \left(\sum_j w_{ij} h_j + b_i \right) \quad (6)$$

où $\sigma(x)$ est la fonction sigmoïde.

6.5 Formalisme

vraisemblance du modèle

Posons

$$\mathcal{L}(\theta | v) = p(v | \theta) = \frac{1}{Z} \sum_h e^{-E(v, h)} \quad (7)$$

La log vraisemblance :

$$\ln \mathcal{L}(\theta | v) = \ln p(v | \theta) = \ln \frac{1}{Z} \sum_h e^{-E(v, h)} = \ln \sum_h e^{-E(v, h)} - \ln \sum_{v, h} e^{-E(v, h)}$$

Gradient de la log vraisemblance :

$$\begin{aligned} \frac{\partial \ln \mathcal{L}(\theta | v)}{\partial \theta} &= \frac{\partial}{\partial \theta} \left(\ln \sum_h e^{-E(v, h)} \right) - \frac{\partial}{\partial \theta} \left(\ln \sum_{v, h} e^{-E(v, h)} \right) \\ &= - \frac{1}{\sum_h e^{-E(v, h)}} \sum_h e^{-E(v, h)} \frac{\partial E(v, h)}{\partial \theta} + \frac{1}{\sum_{v, h} e^{-E(v, h)}} \sum_{v, h} e^{-E(v, h)} \frac{\partial E(v, h)}{\partial \theta} \end{aligned}$$

et en utilisant l'egalit :

$$p(h|v) = \frac{p(v, h)}{p(v)} = \frac{\frac{1}{Z} e^{-E(v, h)}}{\frac{1}{Z} \sum_h e^{-E(v, h)}} = \frac{e^{-E(v, h)}}{\sum_h e^{-E(v, h)}}$$

on a :

$$\frac{\partial \ln \mathcal{L}(\theta|v)}{\partial \theta} = - \sum_h p(h|v) \frac{\partial E(v, h)}{\partial \theta} + \sum_{v, h} p(v, h) \frac{\partial E(v, h)}{\partial \theta} \quad (8)$$

Le passage un lot S de donnes (batch) de taille N donne :

$$\frac{1}{N} \sum_{v \in S} \frac{\partial \ln \mathcal{L}(\theta|v)}{\partial \theta} = \frac{1}{N} \sum_{v \in S} \left[-\mathbb{E}_{p(h|v)} \left[\frac{\partial E(v, h)}{\partial \theta} \right] + \mathbb{E}_{p(h, v)} \left[\frac{\partial E(v, h)}{\partial \theta} \right] \right] \quad (9)$$

$$\frac{1}{N} \sum_{v \in S} \frac{\partial \ln \mathcal{L}(\theta|v)}{\partial \theta} = \frac{1}{N} \sum_{v \in S} \left[-\mathbb{E}_{\text{donnes}} \left[\frac{\partial E(v, h)}{\partial \theta} \right] + \mathbb{E}_{\text{modle}} \left[\frac{\partial E(v, h)}{\partial \theta} \right] \right] \quad (10)$$

Le modele : $\theta = (W, b, c)$

On remplace $\theta = W$ dans le premier terme de l'equation (8) et on a :

$$- \sum_h p(h | v) \frac{\partial E(v, h)}{\partial w_{ij}} = \sum_h p(h | v) h_i v_j$$

Par indpendance des variables conditionnelles $(h_i | v)$ on a : $p(h | v) = \prod_{k=1}^n p(h_k | v)$ et il s'ensuit que ;

$$\begin{aligned} - \sum_h p(h | v) \frac{\partial E(v, h)}{\partial w_{ij}} &= \sum_h \prod_{k=1}^n p(h_k | v) h_i v_j = \sum_{h_i} \sum_{h_{-i}} p(h_i | v) p(h_{-i} | v) h_i v_j \\ &= \sum_h p(h | v) \frac{\partial E(v, h)}{\partial w_{ij}} = \sum_{h_i} p(h_i | v) h_i v_j \underbrace{\sum_{h_{-i}} p(h_{-i} | v)}_{=1} \\ &= \sum_h p(h | v) \frac{\partial E(v, h)}{\partial w_{ij}} = p(H_i = 1 | v) v_j = \sigma \left(\sum_{j=1}^m w_{ij} v_j + c_i \right) v_j \end{aligned}$$

On reprend l'equation (8) et on a :

$$\begin{aligned} \frac{\partial \ln \mathcal{L}(\theta|v)}{\partial w_{ij}} &= - \sum_h p(h|v) \frac{\partial E(v, h)}{\partial w_{ij}} + \sum_{v, h} p(v, h) \frac{\partial E(v, h)}{\partial w_{ij}} \\ &= \sum_h p(h|v) h_i v_j - \sum_v p(v) \sum_h p(h|v) h_i v_j \\ &= p(h_i = 1|v) v_j - \sum_v p(v) p(h_i = 1|v) v_j. \end{aligned} \quad (11)$$

Moyenne du gradient de la log vraisemblance sur un ensemble d'apprentissage

$S = \{\mathbf{v}_1, \dots, \mathbf{v}_N\}$

Les notations suivantes sont souvent utilisées :

$$\begin{aligned} \frac{1}{N} \sum_{v \in S} \frac{\partial \ln \mathcal{L}(\boldsymbol{\theta} | v)}{\partial w_{ij}} &= \frac{1}{N} \sum_{v \in S} \left[-\mathbb{E}_{p(h|v)} \left[\frac{\partial E(v, h)}{\partial w_{ij}} \right] + \mathbb{E}_{p(h,v)} \left[\frac{\partial E(v, h)}{\partial w_{ij}} \right] \right] \\ &= \frac{1}{N} \sum_{v \in S} [\mathbb{E}_{p(h|v)}[v_i h_j] - \mathbb{E}_{p(h,v)}[v_i h_j]] \quad (5) \end{aligned} \quad (12)$$

Par analogie on a ;

$$\frac{\partial \ln \mathcal{L}(\boldsymbol{\theta} | v)}{\partial b_i} = v_j - \sum_v p(v) v_j \quad (13)$$

et

$$\frac{\partial \ln \mathcal{L}(\boldsymbol{\theta} | v)}{\partial c_i} = p(H_i = 1 | v) - \sum_v p(v) p(H_i = 1 | v) \quad (14)$$

Approximation de la log vraisemblance:

L'idée ici est d'éviter un calcul trop complexe qui nécessiterait de parcourir toutes les valeurs possibles des variables visibles (complexité exponentielle). Ce problème apparaît notamment lorsqu'on calcule le gradient du log-vraisemblance, en particulier pour certains termes spécifiques des équations (5), (6) et (7). Les embeddings récupérés en sortie de l'encodeur de Blip sont de tailles $(577)(768) = 443136$, donc il faudrait valuer 2^{443136} possibilités. C'est juste impossible car même le supercalculateur le plus puissant de notre époque ne finirait ce calcul avant des milliers d'années.

Solution possible

Plutôt que d'effectuer ces sommes directement, on peut estimer ces valeurs en prenant des échantillons générés à partir du modèle lui-même. Pour obtenir ces échantillons, on utilise une méthode appelée échantillonnage de Gibbs, qui repose sur une chaîne de Markov. Cette chaîne doit tourner assez longtemps pour atteindre un état stable.

Cependant, même avec cette méthode, le coût de calcul reste trop élevé pour un apprentissage efficace des RBM. C'est pourquoi des approximations supplémentaires sont souvent utilisées, comme la CD-k (k-Contrastive Divergence).

$$CD_k(\boldsymbol{\theta}, \mathbf{v}^{(0)}) = - \sum_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}^{(0)}) \frac{\partial E(\mathbf{v}^{(0)}, \mathbf{h})}{\partial \boldsymbol{\theta}} + \sum_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}^{(k)}) \frac{\partial E(\mathbf{v}^{(k)}, \mathbf{h})}{\partial \boldsymbol{\theta}}.$$

La divergence contrastive (CD) simplifie l'entraînement des RBM en limitant le nombre d'itérations de l'échantillonnage. La chaîne est initialisée avec un exemple d'entraînement $\mathbf{v}(0)$ et génère un échantillon $\mathbf{v}(k)$ après k tapes. À chaque tape, on commence par échantillonner $\mathbf{h}(t)$ à partir de $p(\mathbf{h} | \mathbf{v}(t))$, puis $\mathbf{v}(t+1)$ à partir de $p(\mathbf{v} | \mathbf{h}(t))$. Cette approximation permet d'accélérer le calcul du gradient et l'optimisation du modèle.

Algorithm 1 k-step contrastive divergence

Input: RBM $(V_1, \dots, V_m, H_1, \dots, H_n)$, training batch S

Output: gradient approximation Δw_{ij} , Δb_j , and Δc_i for $i = 1, \dots, n$, $j = 1, \dots, m$

Function k-step contrastive divergence($V_1, \dots, V_m, H_1, \dots, H_n, S$):

```
  for  $i = 1, \dots, n, j = 1, \dots, m$  do
    |  $\Delta w_{ij} = 0, \Delta b_j = 0, \Delta c_i = 0$ 
  end
  for each  $v \in S$  do
    |  $v^{(0)} \leftarrow v$ 
    | for  $t = 0, \dots, k - 1$  do
    |   | for  $i = 1, \dots, n$  do
    |   |   | sample  $h_i^t \sim p(h_i | v^{(t)})$ 
    |   | end
    |   | for  $j = 1, \dots, m$  do
    |   |   | sample  $v_j^{(t+1)} \sim p(v_j | h^{(t)})$ 
    |   | end
    | end
    | for  $i = 1, \dots, n, j = 1, \dots, m$  do
    |   |  $\Delta w_{ij} \leftarrow \Delta w_{ij} + p(H_i = 1 | v^{(0)}) \cdot v_j - p(H_i = 1 | v^{(k)}) \cdot v_j$ 
    |   |  $\Delta b_j \leftarrow \Delta b_j + v_j^{(0)} - v_j^{(k)}$ 
    |   |  $\Delta c_i \leftarrow \Delta c_i + p(H_i = 1 | v^{(0)}) - p(H_i = 1 | v^{(k)})$ 
    | end
  end
end
```

Nous avons introduire la RBM binaire car cest de l que part toutes les variantes existantes de la Machine de Boltzmann Restreinte. Les embeddings obtenues en sortie de lencodeur dimage du modle Blip-Image-Captioning tant des valeurs continues (pas que des 0 et 1), nous allons travailler avec la RBM gaussienne.

7 Quest-ce quune RBM Gaussienne-Binaire?

Cest une **extension** de la RBM binaire qui permet de **traiter des donnes continues** (comme des embeddings, des signaux audio, ou des images en niveaux de gris).

Diffrence cl :

- Dans une RBM **binaire**, les neurones visibles prennent **0 ou 1**.
- Dans une RBM **gaussienne**, les neurones visibles prennent **des valeurs continues**.

Pourquoi cest utile ? Parce que dans ce sujet les embeddings sont **des vecteurs de nombres rels**, donc une RBM gaussienne serait plus adapte quune RBM binaire.

8 Fonctionnement dune RBM-GB

Elle garde la **mme structure** quune RBM classique, mais avec une **modification sur les tats des neurones visibles**.

8.1 Activation des neurones caches (identique la RBM Binaire-Binaire)

Chaque neurone cach h_j reçoit une somme pondérée des entrées v_i et applique une **sigmode** :

$$P(h_j = 1|v) = \sigma \left(\sum_i w_{ij} v_i + b_j \right)$$

Donc les **neurones caches restent binaires (0 ou 1)**, comme dans une RBM classique.

8.2 Activation des neurones visibles (modifié pour les valeurs continues)

Dans une RBM binaire, on n'applique pas une sigmode aux **neurones visibles**. Mais ici, on considère qu'ils suivent une **distribution gaussienne** :

$$v_i \sim \mathcal{N} \left(c_i + \sum_j w_{ij} h_j, \sigma^2 \right)$$

(σ est un hyperparamètre.)

Chaque v_i est un **nombre réel** tiré d'une **distribution normale (gaussienne)** centrée autour de $\sum_j w_{ij} h_j + c_i$, avec une variance σ^2 .

Ce que ça signifie :

- Au lieu d'être **0 ou 1**, les neurones visibles prennent **des valeurs réelles**.
- On peut interpréter ça comme une **reconstruction bruitée, décompressée ou raffinée** des embeddings.

9 Entraînement d'une RBM gaussienne

L'apprentissage suit la même logique qu'une RBM binaire :

1. On passe une **donnée réelle** dans la RBM.
2. On **calcule les activations cachées** avec la sigmode.
3. On **reconstruit les neurones visibles** avec une distribution gaussienne.
4. On ajuste les poids avec le **Contraste de Divergence (CD-k)**.

Petite différence mais importante : Comme les neurones visibles sont continus, la mise à jour des poids prend en compte la variance σ^2 , ce qui change un peu les équations d'apprentissage.

(cette partie est complétée par les nouvelles équations et l'algorithme proprement dit de la RBM gaussienne.)

10 Une fois entrainé, comment on utilise la RBM gaussienne ?

Après l'entraînement, on peut utiliser la RBM gaussienne pour :

- **Transformer des embeddings** : on donne un **vecteur d'embeddings** en entrée, et la RBM génère une **version transformée** des embeddings à partir des activations cachées.
- **Générer de nouvelles données** : on initialise un vecteur au hasard et on le fait passer plusieurs fois dans la RBM, ce qui produit des données similaires à celles du dataset d'entraînement.
- **Extraire des features** : on ne garde que les activations cachées h comme **représentation compacte** des données visibles.

C'est le premier cas d'utilisation qui fait l'objet de ce sujet.

11 Limites d'une RBM gaussienne

- Peut être difficile d'entraîner car la variance doit être bien réglée.
- Complexité computationnelle (parallélisation recommandée)
- Moins populaire car souvent remplacé par des autoencodeurs variationnels (VAE).
- Dans le cas de données complexes (distributions complexes), la RBM Gaussienne peut s'avérer impuissante car elle suppose que les données suivent une gaussienne.

References

- [1] Noam Shazeer, Ashish Vaswani et al. Attention is all you need. <https://arxiv.org/abs/1706.03762>, 2017.
- [2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [3] Yul Choi et al. Spatialvlm: Integrating spatial reasoning with vision-language models. *arXiv preprint arXiv:2105.08752*, 2021.
- [4] Simon Haykin. *Neural Networks: A Comprehensive Foundation*. Prentice Hall, 1998.
- [5] Jialei Hou et al. Vila: Vision-and-language pre-training with large-scale datasets. *arXiv preprint arXiv:2103.12824*, 2021.
- [6] Haesu Hwang et al. Vlm-rl: Visual-language models for reinforcement learning. *arXiv preprint arXiv:2004.06592*, 2020.
- [7] Wonjae Kim et al. Vilt: Vision-and-language transformer without convolution or region supervision. *arXiv preprint arXiv:2102.03334*, 2021.

- [8] Junnan Li et al. Blip: Bootstrapping language-image pre-training. *arXiv preprint arXiv:2201.12086*, 2022.
- [9] Luowei Li et al. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- [10] Wei Liu et al. Xmodel-vlm: Efficient cross-modal representation learning. *arXiv preprint arXiv:2104.06923*, 2021.
- [11] Alec Radford et al. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021.
- [12] Paul Smolensky. Information processing in dynamical systems: Foundations of harmony theory. *Parallel distributed processing: Explorations in the microstructure of cognition*, 1:194–281, 1986.
- [13] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490*, 2019.
- [14] Yao-Hung Hubert Tsai et al. Multimodal transformer for unaligned multimodal language sequences. *arXiv preprint arXiv:1906.00295*, 2019.