# 確率分布を用いた画像テキストデータの埋め込みと検索

Embedding and retrieval of images and text data using probability distribution

濱 健太 *1	松原 崇 *1	上原 邦昭 *1
Kenta Hama	Takashi Matsubara	Kuniaki Uehara

\*1 神戸大学 大学院システム情報学研究科計算科学専攻 Graduate School of System Infomatics, Kobe University

Multimodal data including images, sounds, texts is accumulated on the Internet. We can expect general-purpose data representation to perform tasks such as data discrimination, generation, and retrieval on various modalities datasets. The key idea for acquiring the representation is embedding a point from a data space of each modality in a point of common space. However, if data is embedded in a point, it becomes difficult to interpret the ambiguity of the data's meaning and the inclusive relation among the data. Of course, representation of data point does not necessarily need to be a point. In this study, we embed image and text into a normal distribution in a common space. This improves the performance of image retrieval.

#### 1. はじめに

マルチモーダルデータの蓄積が進み、 モダリティの異なる データ間で,相互の検索や変換を行い,情報の理解,要約をサ ポートするシステムの研究が盛んに行われている [Wang 16]. 複数のモダリティの異なるデータを扱う上では、各モダリティ を横断するデータの分類, 生成, 検索といったタスクに転用が 可能な、汎用的なデータ表現を獲得することが重要である. そ の表現獲得の核となるアイデアとして、各モダリティのデータ からそれら共通の表現空間上の点への埋め込みがある. 埋め込 みとは,あるデータ x をある関数 f によって,データの意味を 表現する空間上のベクトル f(x) に写像することである. さら に、埋め込まれた空間上で類似度を定義すれば、マルチモーダ ルデータ間の意味の近さの測定が可能になる.この空間上では 意味の近いデータ同士の類似度が高いことが望ましいが, デー タ間の類似度を正しく定義されているか評価するのは困難であ る. そのため、埋め込まれた空間の評価に、共通空間上でのデー タ間の検索の精度を利用することが提案されている [Kiros 14]. これは、検索精度の良い共通空間は、意味の近いデータを、適切 な類似度の点に埋め込んでいるという考え方に基づいている.

このようなマルチモーダル表現の獲得を目的として、現在、 盛んに研究されている対象として、画像とテキストデータ間の 検索が挙げられる.画像、テキストデータはインターネット上 から容易に収集でき、画像処理や自然言語処理の分野における 研究結果も多く存在するため、研究対象として選び易い.

画像とテキストデータ間の検索では、まず各データの特 徴抽出を行い、その特徴量を何らかの基準で共通空間上の ベクトルに変換する. 例えば、画像の特徴抽出には ImageNet [Deng 09] で学習済みの CNN (Convolutional Neural Network) の中間層の出力を用いるものや、テキストの特徴抽 出には word2vec [Mikolov 13] や LSTM (Long Short-Term Memory)の出力が用いられる. そして、これらの特徴量を rank-loss と呼ばれるマージンを用いた関数を使い、クエリー に対する望ましい出力ほど類似度が高くなるように埋め込む ことが多い [Kiros 14]. しかし、画像とテキストデータを共通 空間上の点へ埋め込むと、入力データの持つ意味の曖昧さや、 データ間の意味の包含関係といった概念の解釈が困難になる. そのため,データが埋め込まれるべき箇所が複数あるような状 況で問題が発生しやすい.

例えば、単語の埋め込みを考える場合、"rock"という単語は "stone"、"river"などの自然を意味する単語と近い関係にある のと同時に、"pop"、"jazz"などの音楽を意味する単語との意味 も近いはずである.もしもデータを点へ埋め込む場合、"rock" という単語を挟んで、自然を意味する単語全体と音楽を意味す る単語全体が近くの点へ埋め込まれてしまうという問題が生じ る.しかし、データの持つ意味の曖昧さを表現できる埋め込み 手法があれば、"rock"という単語の曖昧さを上手く捉え、周囲 の単語に対する影響を軽減することが期待される.

本研究では、テキストと画像の両方を点ではなく、確率分布 ヘ埋め込む方法を提案する.テキストと画像を確率分布へと埋 め込むことができれば、テキストと画像の持つ意味の曖昧さや、 マルチモーダルデータ間の意味の包含関係や共通部分といった 概念を考えることも可能となる.本研究ではこの確率分布への 埋め込みを行う提案手法を、検索タスクの state-of-the-art で ある VSE++ [Faghri 17] と組み合わせて比較し、類似度の定 義として望ましいことを示す.

#### 2. 関連研究

#### 2.1 VSE (Visual Semantic Embeddings)

ニューラルネットを用いて得られた特徴量から,共通空間 上の点へ埋め込むモデルとして,VSE [Kiros 14] について説 明する.VSE は,画像から説明文を生成するタスク (Image Captioning) を行うためのモデルである.画像の特徴抽出には ImageNet で学習済みの AlexNet [Krizhevsky 12],テキスト の特徴抽出には LSTM を用いて,得られた特徴量は線形変換 によって共通空間上のベクトルへ変換される.

まず, 画像データ x, 画像の特徴量 h, テキストデータ y, デ キストの特徴量 t とする. 画像の特徴量は CNN の最終層から 1 層前の出力で, テキストの特徴量は LSTM の中間層の出力 とし, h = CNN(x), t = LSTM と表すことにする. 次に, こ れらの特徴量を線形変換し, 共通の表現空間に埋め込む. この 埋め込みに用いる行列を各データに対して  $M_x$ ,  $M_y$  とすると, 埋め込まれた共通空間上のベクトルは  $z_h = M_x h$ ,  $z_t = M_y t$ と表される. もし, これらのデータ x, y が関連するデータであ

連絡先: 濱健太, 神戸大学大学院システム情報学研究科計算科 学専攻, hamaken@ai.cs.kobe-u.ac.jp

れば, 共通空間上のベクトル *z<sub>h</sub>*, *z<sub>t</sub>* の類似度も高くなることが 望ましい. VSE では類似度の定義に cos 類似度を用いている. cos 類似度は

$$sim(x,y) = \frac{z_h \cdot z_t}{\|z_h\|_2 \|z_t\|_2}$$

ただし, $z_h = M_x \text{CNN}(x)$ , $z_t = M_y \text{LSTM}(y)$ と定義される. 関連するデータ対 $z_h, z_t$ に関して,  $\sin(x, y)$ が大きくなるように埋め込むために, VSE では rank-loss と呼ばれる以下の式を用いている.

$$r(x,y) = \sum_{\hat{y}} \max\{0, \alpha - \sin(x,y) + \sin(x,\hat{y})\}$$
$$+ \sum_{\hat{x}} \max\{0, \alpha - \sin(x,y) + \sin(\hat{x},y)\}$$

 $\hat{y}, \hat{x}$ はそれぞれ x, yに関連しないデータを表すものとする.  $\alpha$ はマージンと呼ばれるハイパーパラメータであり,関連するデータ対の類似度と,関連しないデータ対の類似度の差をどの程度大きくするかを調節する役割を担っている. データ全体における損失関数は rank-loss を用いて,

$$L = \frac{1}{N} \sum_{n=1}^{N} r(x_n, y_n)$$

と定義する.ただし,Nはデータ数を表すものとする.定義した損失関数を最小化するように学習を行えば,共通空間への埋め込みが実現される.

#### 2.2 VSE++

VSE++ (Improved Visual Semantic Embeddings) [Faghri 17] は VSE の改良として,提案されたものである. VSE の rank-loss を以下のように変更している.

$$r(x,y) = \max_{\hat{y}} \max\{0, \alpha - \sin(x,y) + \sin(x,\hat{y})\}$$
$$+ \max_{\hat{x}} \max\{0, \alpha - \sin(x,y) + \sin(\hat{x},y)\}$$

VSE では関連しないデータ点全てにおいての和をとっていた 箇所が, 各データ点における損失の中でも最大の値となるデー タ点のみを全体の損失に含めている. シンプルな変更であるが, VGG [Simonyan 14], ResNet [He 16] など画像の特徴抽出に 使うモデルによらず精度が向上し, 別の埋め込み手法との併用 も可能である. 本研究ではこの VSE++を基本モデルとし, 確 率分布への埋め込みが可能となる改良を加える.

### 3. 提案手法

本研究では、各入力画像、テキストデータを多変量正規分布 に変換する.その際、計算を簡単にするために、正規分布の各 次元に独立性を仮定し、分散共分散行列は対角成分のみ持つこ とにしている.提案手法のモデル全体は図1のようになって いる.

まず入力データを VSE++と同様に、画像は ImageNet で事前学習済みの CNN (VGG, ResNet) に入力し、テキ ストは LSTM に入力する. なお本研究における実装には GRU [Cho 14] を用いている. 各ニューラルネットワーク を通して得られた特徴量は、線形変換によって正規分布の パラメータである平均と分散に変換される. 埋め込みは  $M_{x_{\mu}}, M_{x_{\sigma}}, M_{y_{\mu}}, M_{y_{\sigma}}$ の4つの行列を用いて行われる. これら の行列を用いて、図1の入力画像データとテキストデータの平 均と分散は、 $\mu_i = M_{x_{\mu}}i, \sigma_i = \text{diag}(M_{x_{\sigma}}i), \mu_t = M_{y_{\mu}}t, \sigma_t = \text{diag}(M_{y_{\sigma}}t)$ と表される.よって、画像 x の特徴量 i とテキスト y の特徴量 t の埋め込みは、 $z_i = \mathcal{N}(\mu_i, \sigma_i), z_t = \mathcal{N}(\mu_t, \sigma_t)$ という正規分布の形で表現される.以上の操作と VSE++は、CNN、LSTM を用いて特徴抽出を行うところまでは同じで、その後、特徴量を1つの点ではなく、2つの平均と、分散を表現する点へ変換するように埋め込み行列を2つずつ用意する点で 異なっている.



 $\boxtimes$  1: The whole model of the proposed method.

この埋め込みは、画像、テキストの特徴量を n 次元実数空間 上の点  $x \in R^n$  ではなく、多変量正規分布の平均  $\mu \in R^n$ ,分 散共分散行列  $\Sigma \in R^{(n,n)}$ (ただし  $\Sigma = \text{diag}(\sigma), \sigma \in R^n$ ) へ と変換すれば、関数空間上の点  $f \in \{N(\mu, \Sigma) | \mu \in R^n, \Sigma = \text{diag}(\sigma), \sigma \in R^n\}$  への埋め込みを行なっていると解釈できる. よって、VSE++の共通空間である  $R^n$  と空間そのものが異な るため、cos 類似度をそのまま利用することはできない. そこ で、確率分布間の類似度を定義するため、確率分布間の差異を 計算する際に用いられることの多い、JS ダイバージェンスを 利用する.

JS ダイバージェンスは, KL ダイバージェンスを用いて定義 される. 二つの確率分布を p, q とすると, KL ダイバージェン スは  $D_{KL}(p||q) = \int p \log \frac{p}{q} dx$  で定義される. p,q は正規分布 なので,  $D_{KL}(p||q)$  は代数計算によって求めることができ, さ らに各次元の独立性を仮定しているので, 計算も容易である. JS ダイバージェンスは KL ダイバージェンスを用いて

$$D_{JS}(p||q) = \frac{1}{2}(D_{KL}(p||q) + D_{KL}(q||p))$$

と定義される. ここで類似度は符号を反転し,  $sim(x, y) = -D_{JS}(p||q)$ と定義する. もちろん, 類似度の定義に KL ダイ バージェンスを用いることも考えられる. しかし, 事前実験に より, 類似度に KL ダイバージェンスより JS ダイバージェン スを用いた方が精度が良かったため, 本研究では JS ダイバー ジェンスによる埋め込みを行っている.

モデルの学習は二段階で行う.まず第一段階として,画像 データ,テキストデータの平均に対して,通常のVSE++と同 様に, cos 類似度を用いて埋め込みを行う.その後,第二段階と して,各データの分散を学習する.この際,一段階目で学習し たテキストデータの平均を固定して学習を行う.これは,事前 実験によりデータの平均を固定すると,固定せず学習した場合 より精度の向上が見られたためである.また,第二段階の分散 の学習では先ほど定義した JS ダイバージェンスを用いた類似 度を使う.

	Caption Retrieval				Image Retrieval					
Model	R@1	R@5	R@10	Med r	Mean r	R@1	R@5	R@10	Med r	Mean r
VSE++(VGG) Ours(VGG)	52.0 <b>52.6</b>	82.0 <b>82.3</b>	89.9 <b>90.5</b>	$\begin{array}{c} 1.0\\ 1.0\end{array}$	5.3 <b>5.2</b>	39.6 <b>40.5</b>	74.5 <b>75.0</b>	85.7 <b>86.2</b>	$2.0 \\ 2.0$	9.6 <b>8.5</b>
VSE++(ResNet) Ours(ResNet)	59.2 <b>59.4</b>	86.1 <b>86.7</b>	93.2 <b>93.8</b>	$\begin{array}{c} 1.0\\ 1.0\end{array}$	4.0 <b>3.8</b>	43.8 <b>44.9</b>	78.0 <b>78.9</b>	88.2 88.9	$2.0 \\ 2.0$	7.6 <b>6.7</b>
VSE++(VGG, finetune) Ours(VGG, finetune)	57.6 <b>58.4</b>	86.0 <b>86.5</b>	93.0 <b>93.6</b>	$\begin{array}{c} 1.0\\ 1.0\end{array}$	$\begin{array}{c} 4.1\\ 4.1\end{array}$	45.2 <b>45.4</b>	79.1 <b>79.3</b>	88.9 88.9	$\begin{array}{c} 2.0 \\ 2.0 \end{array}$	8.3 <b>7.5</b>
VSE++(ResNet, finetune) Ours(ResNet, finetune)	$\begin{array}{c} 64.3\\ 64.3\end{array}$	89.5 <b>90</b>	95.9 <b>96.3</b>	$\begin{array}{c} 1.0\\ 1.0\end{array}$	$\begin{array}{c} 3.1\\ 3.1\end{array}$	50.0 <b>50.5</b>	83.3 83.3	91.4 91.4	1.6 <b>1.4</b>	6.0 <b>5.4</b>

表 1: Effect of feature and fine-tuning.

# 4. 実験及び結果

学習データセットとして, Image Caption Retrieval で一般 に用いられる, Microsoft COCO データセット [Lin 14](以下 MS COCO) を使用し, 画像テキスト間の双方向の検索を評価 する. MS COCO は, 物体検知, セグメンテーション, キャプ ショニングのための大規模な画像とテキストのデータセットで ある. Image Caption Retrieval では, キャプショニング用の データを用いる. キャプショニング用のデータは, 1枚の画像 に対して, その説明を行うテキスト (キャプション) が5つ与 えられている. データの分割は, VSE++の論文と同様に, 訓 練データに 113,287枚, 検証データに 5000枚, 評価データに 5000枚の画像を使用している. 画像データは, CNN への入力 時に画像全体からランダムに 224 × 224 の大きさで切り出し ている.

パラメータの最適化アルゴリズムには、Adam [Kingma 14] を用いている. ハイパーパラメータは、一段階目のエポック数 を 30、学習率は最初の 15 エポックを  $2.0 \times 10^{-4}$ ,残りの 15 エ ポックを  $2.0 \times 10^{-5}$  に設定している. 二段階目の学習はエポッ ク数を 30、学習率は最初の 15 エポックを  $2.0 \times 10^{-5}$ ,残りの 15 エポックを  $2.0 \times 10^{-6}$  に設定している. 学習率が一段階目 に対して 0.1 倍されているのは、画像の平均の固定を行なって いないため、一段階目で学習された平均の局所解を抜け出して、 別の局所解へ向かうことを防ぐためである.

また,本研究では CNN の fine-tuning と提案手法を組み合 わせた実験も行なっている. CNN の fine-tuning は,一段階目 が終了した後に CNN のパラメータの固定をはずし,エポック 数を 15, 学習率  $2.0 \times 10^{-6}$  に設定している. その後, CNN の パラメータを固定して,先述した二段階目の学習と同様の設定 で分散を学習させている. これらの学習の rank-loss のマージ ンは, fine-tuning を行わない場合は 0.2 を,行う場合は 0.15 を 使用している. 検索結果の評価指標には一般に Image Caption Retrieval で用いられる R@k, Med r, Mean r を用いている. 評価データ 5000 枚は 1000 枚ずつに対して検索を行い,5 つの 評価値の平均を評価結果としている.

実験では、4 つの条件設定において既存手法である VSE++ と提案手法である正規分布への埋め込みを比較する.4 つの条 件とは、CNN に VGG か ResNet を用いた場合を、それぞれ fine-tuning ありの場合となしの場合の4 パターンに対応して いる.これらの比較は、VGG や ResNet の出力といった異なる 特徴量を用いた場合にも本手法が適用できること、fine-tuning と組み合わせることが可能であることを示すために行っている. 結果を表1に示す.

表1の上から、VGG, ResNet, VGG + fine-tuning, ResNet + fine-tuning の場合における既存手法(上)と提案手法(下)の比較である. R@1,R@5,R@10,Med r, Mean r 全ての指標において,提案手法の精度が向上している. また画像検索(Image Retrieval)とテキスト検索(Caption Retrieval)を比較すると,評価値の改善は画像検索の mean r が顕著である. 特に, VGG, ResNet, VGG + fine-tuning, ResNet + fine-tuning の各条件で Caption Retrieval の Mean r は平均 0.075 の改善であるの に対して, Image Retrieval の Mean r は平均 0.85 ほど小さくなっている. 以上のことから,提案手法が画像検索側の Mean r に対して,大きく影響を及ぼしていることが分かる.

## 5. 考察

表1より, VGG, ResNet の中間層の出力どちらを用いた場 合においても, 精度の向上が見られていることから, 提案手法 が画像特徴量によらず適用可能であると考えられる.この事実 から, 画像の分類タスクにおける CNN モデルの改良が進んだ 場合, そのモデルに提案手法を適用すれば, 検索精度のさらな る向上が期待できる.

次に、Meanrの精度が向上した理由について考察を行う.図 2は、画像検索において Mean r を大きく改善させた例の一つ である. 図2下の表における Rank は, クエリーとなるテキス トに対して正解となる上の画像の出力順位を示している.上か ら4番目のテキストの Rank を見ると、VSE++は 665 と他と 比較してかなり大きい値を取っている. テキストの内容に注目 すると、このテキストのみ "bus" という単語が含まれていな い. これは、VSE++が LSTM の学習において、テキスト内の "bus"という単語の有無のみで、画像とテキストを類似度の高 い点へ埋め込んでしまっていると考えられる.しかし、分布へ の埋め込みを行う提案手法の Rank は 282 であり, 大きく改善 している.これは図2上のバスの画像の分散を大きくとるか, "bus" という単語の含まれていない, 上から4番目のテキスト の分散を大きくとるように、提案手法が学習したからだと考え られる.このことから、提案手法は各データの持つ意味の曖昧 さを捉えていると考えられる.



	Rank		
Caption	VSE++	Ours	
Passenger bus stationary in traffic on city street.	3	2	
A modern city commuter bus in traffic during the day	2	1	
A bus parked on the side of the road while in traffic	5	6	
a modern train is parked against the sidewalk curb	665	282	
A bus stops at a curb in a busy city street.	6	3	

 $\boxtimes$  2: An example of greatly improving Mean r.

表1の fine-tuning ありとの比較を見ると, 評価値の向上が 小さいものも見られる. これは, fine-tuning することによって, モデルの表現力が大きくなっていることが関係すると考えられ る. ImageNet での事前学習に用いる画像の枚数は 120 万枚ほ どなのに対して, MS COCO の訓練データ数は 11 万枚ほどし かないので, fine-tuning によって, 平均の学習で過学習を起こ している可能性がある. そのため, その後の分散の学習による 精度の向上が困難になったと考えられる. しかし, 画像検索の mean r に関しては大きく改善しているため, 提案手法による 改善が, 誤差の影響によるものとは考えにくい.

# 6. 結論

本研究では VSE++の埋め込みの出力を正規分布の平均と 分散に変更し, rank-loss の類似度を JS ダイバージェンスを用 いて定義するという変更を加えて, 画像データとテキストデー タの共通空間上の点への埋め込みから, 確率分布への埋め込み に拡張する方法を提案した.提案手法は学習を行う際, 通常 の VSE++と同様に学習を行うのではなく, 平均のみを事前に VSE++で学習した上で, 平均を固定して追加で分散の学習を 行うために, 分散を上手く学習することが可能となり, 検索精 度が向上された. この結果から,提案手法が共通空間上で望ま しい類似度を定義していることの傍証が得られた.

本研究の今後の課題について述べる.まず埋め込みに用い る確率分布を変更するといった工夫が可能である.本研究で正 規分布に仮定した各次元の独立性を外すか,多峰性の混合正規 分布に埋め込むことも考えられる.また,分布への埋め込みを 行っていることから,新たな正則化の導入も考えられる.例と して,KLダイバージェンスを用いて分布間の包含関係を表現 すれば,ランダムクロップされた画像の分布が切り出される前 の画像に包含されるといった仮定を,正則化条件として加える ことができる.このような点への埋め込みでは導入できない, 分布への埋め込みのための正則化によって,モデルの汎化能力 が向上することが期待される.

本研究は科研費 (16K12487) の支援,総務省 SCOPE(受付 番号 172107101) の委託を受けて行われた.

# 参考文献

- [Cho 14] Cho, K., Merrienboer, van B., Gülçehre, Ç., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y.: Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation, in *Proceedings of the Conference on Empirical Methods in Natural Language Processing(EMNLP)*, pp. 1724–1734 (2014)
- [Deng 09] Deng, J., Dong, W., Socher, R., Li, L., Li, K., and Li, F.: ImageNet: A large-scale hierarchical image database, in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition(CVPR)*, pp. 248–255 (2009)
- [Faghri 17] Faghri, F., Fleet, D. J., Kiros, R., and Fidler, S.: VSE++: Improved Visual-Semantic Embeddings, arXiv, pp. 1–11 (2017)
- [He 16] He, K., Zhang, X., Ren, S., and Sun, J.: Deep Residual Learning for Image Recognition, in *Proceedings* of the IEEE Conference on Computer Vision and Pattern Recognition(CVPR), pp. 770–778 (2016)
- [Kingma 14] Kingma, D. P. and Ba, J.: Adam: A Method for Stochastic Optimization, arXiv, pp. 1–15 (2014)
- [Kiros 14] Kiros, R., Salakhutdinov, R., and Zemel, R. S.: Unifying Visual-Semantic Embeddings with Multimodal Neural Language Models, arXiv, pp. 1–13 (2014)
- [Krizhevsky 12] Krizhevsky, A., Sutskever, I., and Hinton, G. E.: ImageNet Classification with Deep Convolutional Neural Networks, in *Proceedings of the 26th An*nual Conference on Neural Information Processing Systems(NIPS), pp. 1106–1114 (2012)
- [Lin 14] Lin, T., Maire, M., Belongie, S. J., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L.: Microsoft COCO: Common Objects in Context, in *Pro*ceedings of the 13th European Conference of Computer Vision(ECCV), pp. 740–755 (2014)
- [Mikolov 13] Mikolov, T., Chen, K., Corrado, G., and Dean, J.: Efficient Estimation of Word Representations in Vector Space, arXiv, pp. 1–12 (2013)
- [Simonyan 14] Simonyan, K. and Zisserman, A.: Very Deep Convolutional Networks for Large-Scale Image Recognition, arXiv, pp. 1–14 (2014)
- [Wang 16] Wang, W., Yang, X., Ooi, B. C., Zhang, D., and Zhuang, Y.: Effective deep learning-based multimodal retrieval, Very Large Data Bases Journal(VLDB J), Vol. 25, No. 1, pp. 79–101 (2016)