

Multimodal Deep Learning

"Unify multimodal signals into a single vector space and thereby enable cross-modality signal processing."

Representation Learning

Multimodal Representation

Unimodal Embeddings:

1. Visual Representation: image embeddings
2. Language Representation: Text embeddings

3. Vector arithmetic

for word and image embeddings

Word embeddings

Syntactic / semantic regularities

4. Speech Representation: i-vector approach

Multimodal Representation:

joint embeddings to leverage the complementarity of multimodal data to represent such concepts more accurately

① Unsupervised Training Methods:

deep Boltzmann machines

autoencoders

deep multimodal similarity model (DMSM)

generate fine-grained multimodal embeddings

deep attentional multimodal similarity model (DAMSM)

"measure the similarity between images"

sub-regions and words as an
additional loss function for
text-to-image generation"

② Supervised Training Methods

improve the learning of multimodal representation

Discriminative
factors

(Supervised training)

Intra-modality
generative factors

(Unsupervised training)

③ Methods for zero-shot Learning

"certain representations may require
pairwise data from different modalities

Simultaneously."

④ Transformer-based Methods:

↓
Fusion of Multimodal Signals

"Integrate information extracted from different unimodal data sources into a single compact multimodal representation"

① Simple Operation-based Fusion

② Attention-based Fusion

③ Bilinear Pooling-based Fusion

{ Factorization for Bilinear Pooling

Bilinear Pooling and Attention Mechanisms

Applications

// Image Captioning

// Text-to-Image Generation

① GAN-based Methods

② Generating High-quality images

③ Generating Semantically consistent images

④ Semantic layout control for complex scenes

// Visual Question Answering

// Visual Reasoning