

PLUG AND PLAY AUTOENCODERS FOR CONDITIONAL TEXT GENERATION

FLORIAN MAI †, ♠ NIKOLAOS PAPPAS ♣
IVAN MONTERO ♣ NOAH A. SMITH ♣, ♦ JAMES HENDERSON †

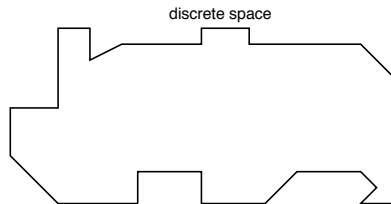
† IDIAP RESEARCH INSTITUTE, ♠ EPFL, SWITZERLAND

♣ UNIVERSITY OF WASHINGTON, SEATTLE, USA

♦ ALLEN INSTITUTE FOR ARTIFICIAL INTELLIGENCE, SEATTLE, USA

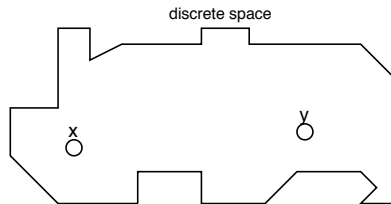


THE PROBLEM WITH CONDITIONAL TEXT GENERATION



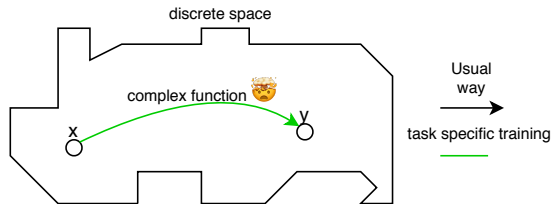
- text lives in a messy, discrete space

THE PROBLEM WITH CONDITIONAL TEXT GENERATION



- text lives in a messy, discrete space
- conditional text generation requires mapping from discrete input to discrete output

THE PROBLEM WITH CONDITIONAL TEXT GENERATION

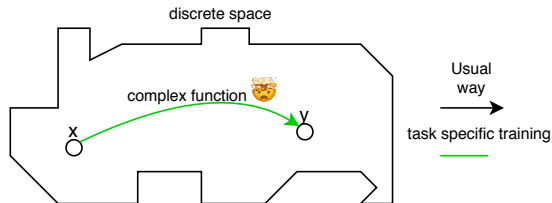


- text lives in a messy, discrete space
- conditional text generation requires mapping from discrete input to discrete output

Usual way:

- learning a complex, task-specific function, which is difficult to train in discrete space

THE PROBLEM WITH CONDITIONAL TEXT GENERATION



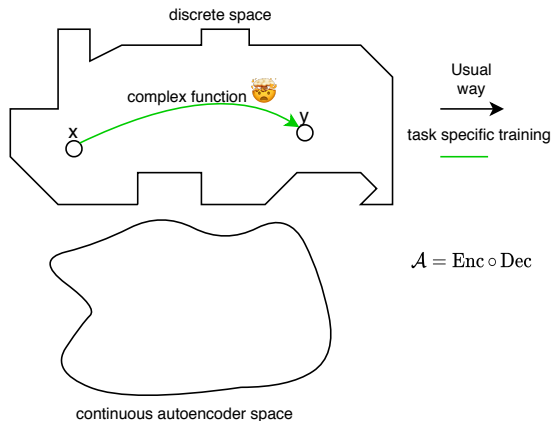
- text lives in a messy, discrete space
- conditional text generation requires mapping from discrete input to discrete output

Usual way:

- learning a complex, task-specific function, which is difficult to train in discrete space

Our way:

THE PROBLEM WITH CONDITIONAL TEXT GENERATION



- text lives in a messy, discrete space
- conditional text generation requires mapping from discrete input to discrete output

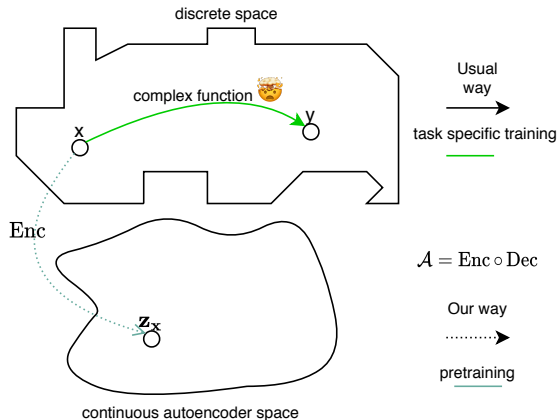
Usual way:

- learning a complex, task-specific function, which is difficult to train in discrete space

Our way:

- obtain a continuous space by training an autoencoder

THE PROBLEM WITH CONDITIONAL TEXT GENERATION



- text lives in a messy, discrete space
- conditional text generation requires mapping from discrete input to discrete output

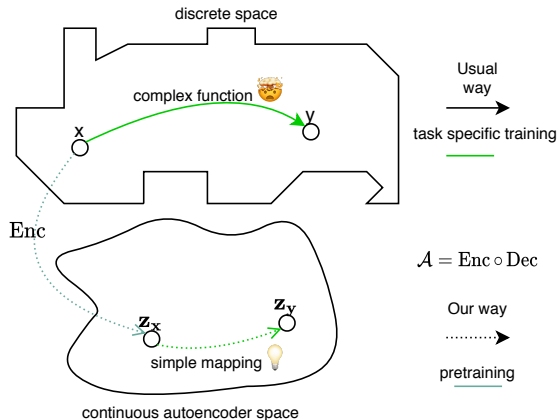
Usual way:

- learning a complex, task-specific function, which is difficult to train in discrete space

Our way:

- obtain a continuous space by training an autoencoder

THE PROBLEM WITH CONDITIONAL TEXT GENERATION



- text lives in a messy, discrete space
- conditional text generation requires mapping from discrete input to discrete output

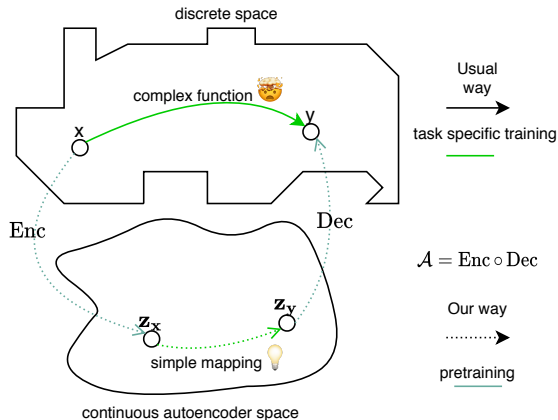
Usual way:

- learning a complex, task-specific function, which is difficult to train in discrete space

Our way:

- obtain a continuous space by training an autoencoder
- reduce task-specific learning to the continuous space

THE PROBLEM WITH CONDITIONAL TEXT GENERATION



- text lives in a messy, discrete space
- conditional text generation requires mapping from discrete input to discrete output

Usual way:

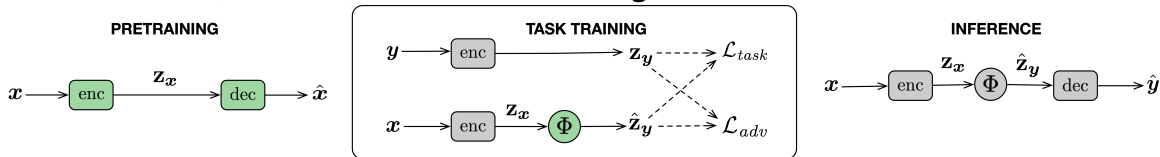
- learning a complex, task-specific function, which is difficult to train in discrete space

Our way:

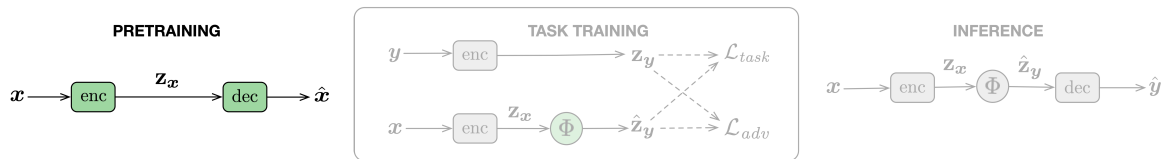
- obtain a continuous space by training an autoencoder
- reduce task-specific learning to the continuous space

FRAMEWORK OVERVIEW

Our framework (*Emb2Emb*) consists of three stages:



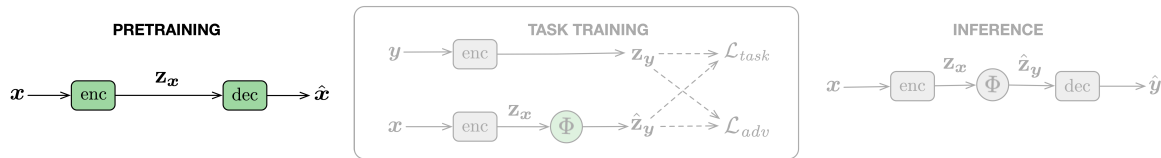
FRAMEWORK OVERVIEW



Pretraining:

- Train a model of the form $\mathcal{A}(x) = \text{Dec}(\text{Enc}(x))$ on corpus of sentences
- Assume a fixed-size continuous embedding $z_x := \text{Enc}(x) \in \mathbb{R}^d$
- Enc and Dec can be any function trained with any objective so long as $\mathcal{A}(x) \approx x$
- training corpus can be any unlabeled corpus \Rightarrow large-scale pretraining?

FRAMEWORK OVERVIEW



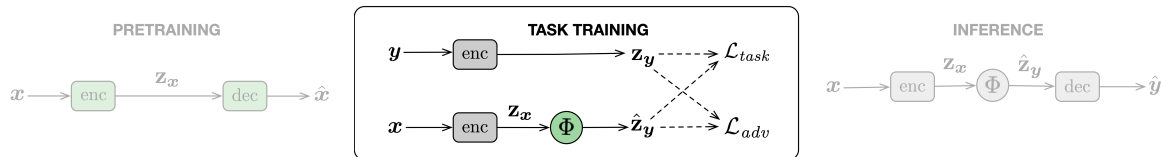
Pretraining:

- Train a model of the form $\mathcal{A}(x) = \text{Dec}(\text{Enc}(x))$ on corpus of sentences
- Assume a fixed-size continuous embedding $z_x := \text{Enc}(x) \in \mathbb{R}^d$
- Enc and Dec can be any function trained with any objective so long as $\mathcal{A}(x) \approx x$
- training corpus can be any unlabeled corpus \Rightarrow large-scale pretraining?

Plug and Play

Our framework is *plug and play* because any autoencoder can be used with it.

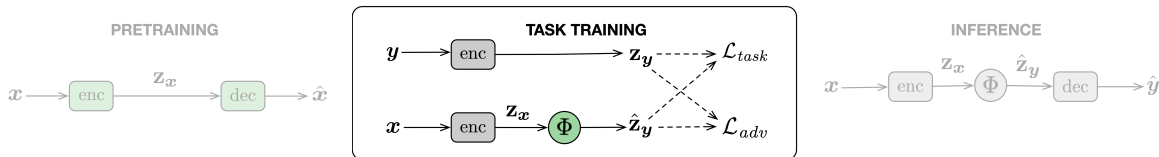
FRAMEWORK OVERVIEW



Task Training:

- Supervised case: $\mathcal{L}_{task}(\hat{z}_y, z_y) = d(\hat{z}_y, z_y)$ where d is a distance function (cosine distance loss in our experiments).

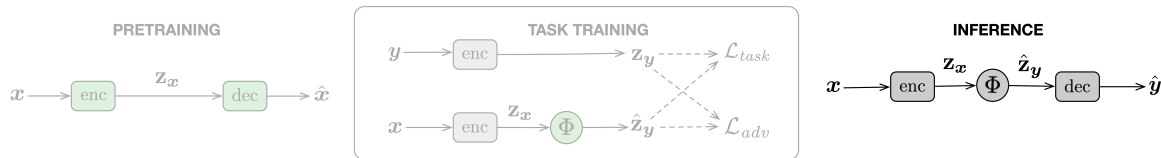
FRAMEWORK OVERVIEW



Task Training:

- Supervised case: $\mathcal{L}_{task}(\hat{\mathbf{z}}_y, \mathbf{z}_y) = d(\hat{\mathbf{z}}_y, \mathbf{z}_y)$ where d is a distance function (cosine distance loss in our experiments).
- Training objective: $\mathcal{L} = \mathcal{L}_{task} + \lambda_{adv} \cdot \boxed{\mathcal{L}_{adv}^?}$

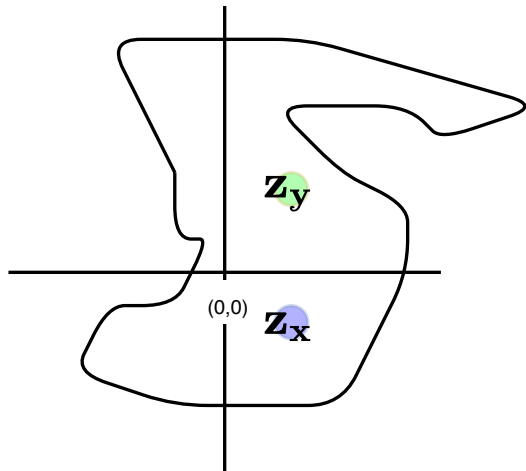
FRAMEWORK OVERVIEW



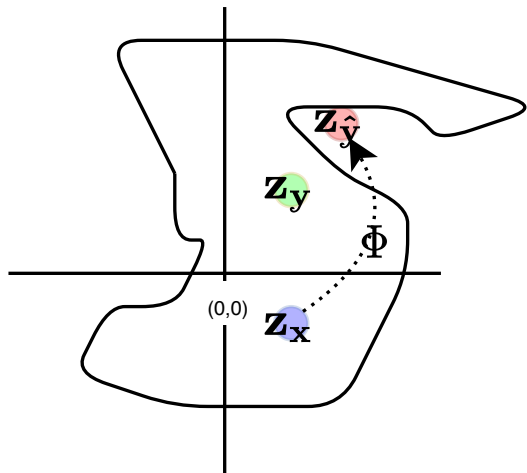
Inference:

- compose inference model as $\text{Enc} \circ \Phi \circ \text{Dec}$
- but: Dec not involved in training. Can it handle outputs of Φ ?
- \Rightarrow yes, if using \mathcal{L}_{adv} .

WHAT CAN HAPPEN WHEN LEARNING IN THE EMBEDDING SPACE?

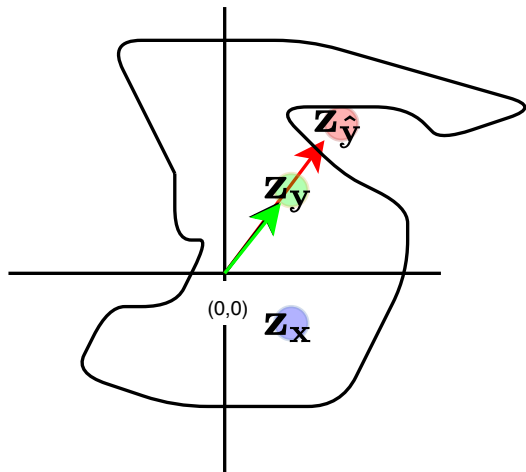


WHAT CAN HAPPEN WHEN LEARNING IN THE EMBEDDING SPACE?



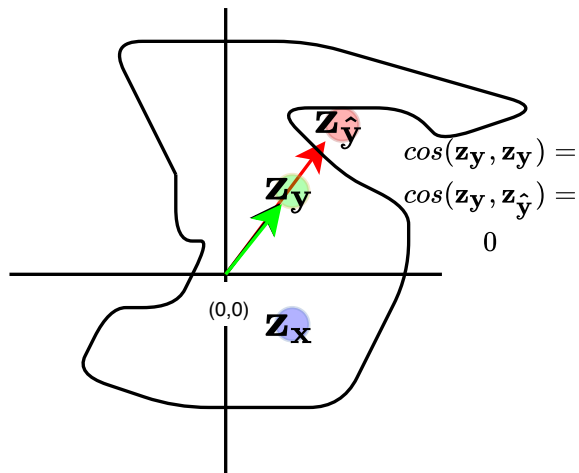
- A prediction may end up **off** the manifold, and by definition, the decoder cannot handle off-manifold data well, but ...

WHAT CAN HAPPEN WHEN LEARNING IN THE EMBEDDING SPACE?



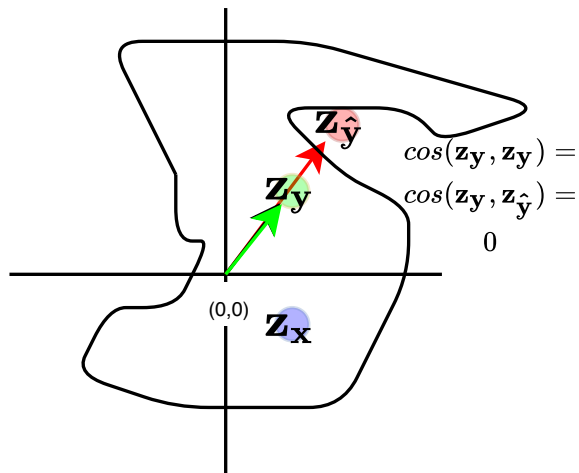
- A prediction may end up **off** the manifold, and by definition, the decoder cannot handle off-manifold data well, but ...
- ... but the predicted embedding may still have the same angle as the true output embedding...

WHAT CAN HAPPEN WHEN LEARNING IN THE EMBEDDING SPACE?



- A prediction may end up **off** the manifold, and by definition, the decoder cannot handle off-manifold data well, but ...
- ... but the predicted embedding may still have the same angle as the true output embedding...
- resulting in zero cosine distance loss despite being off the manifold.

WHAT CAN HAPPEN WHEN LEARNING IN THE EMBEDDING SPACE?



- A prediction may end up **off** the manifold, and by definition, the decoder cannot handle off-manifold data well, but ...
- ... but the predicted embedding may still have the same angle as the true output embedding...
- resulting in zero cosine distance loss despite being off the manifold.
- Similar problems arise for L2 distance - how do we keep the embeddings **on** the manifold?

ADVERSARIAL LOSS TERM

- train a discriminator disc to distinguish between embeddings produced by the encoder and embeddings resulting from the mapping:

$$\max_{\text{disc}} \sum_{i=1}^N \log(\text{disc}(\mathbf{z}_{\tilde{\mathbf{y}}_i})) + \log(1 - \text{disc}(\Phi(\mathbf{z}_{\mathbf{x}_i})))$$

- using the adversarial learning framework, mapping acts as the adversary and tries to fool the discriminator:

$$\mathcal{L}_{adv}(\Phi(\mathbf{z}_{\mathbf{x}_i}); \theta) = -\log(\text{disc}(\Phi(\mathbf{z}_{\mathbf{x}_i}); \theta))$$

- at convergence, the mapping should only produce embeddings that are on the manifold

SUPERVISED STYLE TRANSFER EXPERIMENTS

- *WikiLarge* dataset: transform “normal” English to “simple” English
- parallel sentences (input and output) are available

Model	BLEU (relative imp.)	SARI (relative imp.)
Emb2Emb (no \mathcal{L}_{adv})	15.7 (-)	21.1 (-)
Emb2Emb	34.7 (+121%)	25.4 (+20.4%)

The adversarial loss term \mathcal{L}_{adv} is crucial for embedding-to-embedding training!

SUPERVISED STYLE TRANSFER EXPERIMENTS

- we conducted controlled experiments of models **with a fixed-size bottleneck**
- best Seq2Seq model: best performing variant among fixed-size bottleneck models that are trained end-to-end via token-level cross-entropy loss (like Seq2Seq)

Model	BLEU (relative imp.)	SARI (relative imp.)	Speedup
Best Seq2Seq model	23.3 ($\pm 0\%$)	22.4 ($\pm 0\%$)	-
Emb2Emb	34.7 (+48.9%)	25.4 (+13.4%)	2.2 \times

Training models with a fixed-size bottleneck may be *easier*, *faster*, and *more effective* when training embedding-to-embedding!

UNSUPERVISED TASK TRAINING

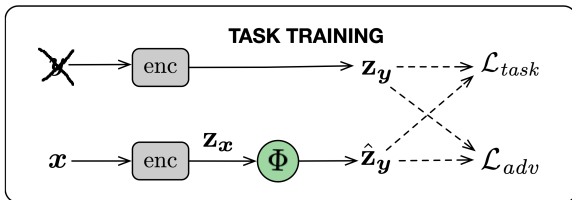
- Fixed-size bottleneck autoencoders are commonly used for unsupervised style transfer

UNSUPERVISED TASK TRAINING

- Fixed-size bottleneck autoencoders are commonly used for unsupervised style transfer
- The goal is to change the **style** of a text, but retain the **content**: e.g., in machine translation, sentence simplification, sentiment transfer

UNSUPERVISED TASK TRAINING

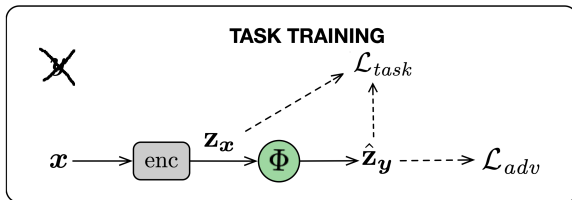
- Fixed-size bottleneck autoencoders are commonly used for unsupervised style transfer
- The goal is to change the **style** of a text, but retain the **content**: e.g., in machine translation, sentence simplification, sentiment transfer



- training objective: $\mathcal{L} = \mathcal{L}_{task} + \lambda_{adv} \cdot \mathcal{L}_{adv}$

UNSUPERVISED TASK TRAINING

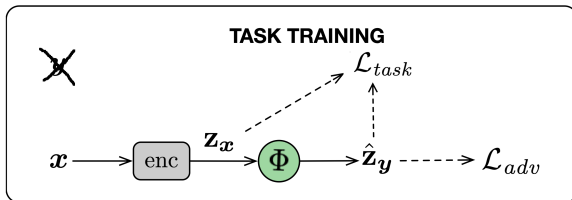
- Fixed-size bottleneck autoencoders are commonly used for unsupervised style transfer
- The goal is to change the **style** of a text, but retain the **content**: e.g., in machine translation, sentence simplification, sentiment transfer



- training objective: $\mathcal{L} = \mathcal{L}_{task} + \lambda_{adv} \cdot \mathcal{L}_{adv}$
- $\mathcal{L}_{task}(\hat{\mathbf{z}}_y, \mathbf{z}_x) = \lambda_{sty} \mathcal{L}_{sty}(\hat{\mathbf{z}}_y) + (1 - \lambda_{sty}) \mathcal{L}_{cont}(\hat{\mathbf{z}}_y, \mathbf{z}_x)$

UNSUPERVISED TASK TRAINING

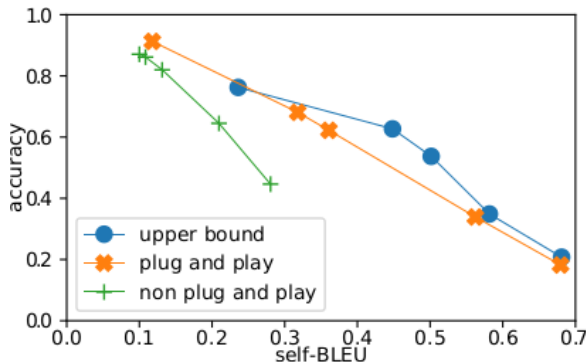
- Fixed-size bottleneck autoencoders are commonly used for unsupervised style transfer
- The goal is to change the **style** of a text, but retain the **content**: e.g., in machine translation, sentence simplification, sentiment transfer



- training objective: $\mathcal{L} = \mathcal{L}_{task} + \lambda_{adv} \cdot \mathcal{L}_{adv}$
- $\mathcal{L}_{task}(\hat{\mathbf{z}}_y, \mathbf{z}_x) = \lambda_{sty} \mathcal{L}_{sty}(\hat{\mathbf{z}}_y) + (1 - \lambda_{sty}) \mathcal{L}_{cont}(\hat{\mathbf{z}}_y, \mathbf{z}_x)$
- we set \mathcal{L}_{cont} to cosine distance, and \mathcal{L}_{sty} to a style classifier's negative log-likelihood of the target class

UNSUPERVISED STYLE TRANSFER EXPERIMENTS

- *Yelp* sentiment transfer dataset: transform reviews with negative sentiment into reviews with positive sentiment (accuracy), but retain content (self-BLEU)
- if we have labels for only 10% of the data, how much better is a plug and play model?



Effect of pretraining

By leveraging autoencoder pretraining on unlabeled data, our plug and play method offers a distinct advantage on unsupervised style transfer!

CONCLUSION

In this talk...

CONCLUSION

In this talk...

- we propose to learn in the embedding space of a pretrained autoencoder, training embedding-to-embedding (Emb2Emb).

CONCLUSION

In this talk...

- we propose to learn in the embedding space of a pretrained autoencoder, training embedding-to-embedding (Emb2Emb).
- we discuss why it's important to keep the predicted embedding on the manifold of the autoencoder, and how to achieve that.

CONCLUSION

In this talk...

- we propose to learn in the embedding space of a pretrained autoencoder, training embedding-to-embedding (Emb2Emb).
- we discuss why it's important to keep the predicted embedding on the manifold of the autoencoder, and how to achieve that.
- we demonstrate that a plug and play method like ours has a distinct advantage on unsupervised style transfer.

CONCLUSION

In this talk...

- we propose to learn in the embedding space of a pretrained autoencoder, training embedding-to-embedding (Emb2Emb).
- we discuss why it's important to keep the predicted embedding on the manifold of the autoencoder, and how to achieve that.
- we demonstrate that a plug and play method like ours has a distinct advantage on unsupervised style transfer.

Additionally, our paper...

- presents an architecture for the mapping Φ that is better than just MLPs.

CONCLUSION

In this talk...

- we propose to learn in the embedding space of a pretrained autoencoder, training embedding-to-embedding (Emb2Emb).
- we discuss why it's important to keep the predicted embedding on the manifold of the autoencoder, and how to achieve that.
- we demonstrate that a plug and play method like ours has a distinct advantage on unsupervised style transfer.

Additionally, our paper...

- presents an architecture for the mapping Φ that is better than just MLPs.
- demonstrates how to further improve the performance on unsupervised style transfer at inference time.

THANK YOU