CS-GY 6923: Lecture 13 Principal Component Analysis, Semantic Embeddings

NYU Tandon School of Engineering, Prof. Christopher Musco

TRANSFER LEARNING

Empirical observation: Features learned when training models like deep neural nets are often useful for problems beyond what the model was trained on.



Very useful in domains like computer vision where we have huge labeled datasets to train deep models on. Approach:

- 1. Download network trained on large image classification dataset (e.g. Imagenet).
- 2. Extract features <u>z</u> for any new image x by running it through the network up until layer before last.
- 3. Use these features for a new problem (e.g., quidditch ball detection), typically using a simpler machine learning algorithm that requires less data (nearest neighbor, logistic regression, etc.).

But what if we don't even have labeled data for a sufficiently related problem?

How to extract features in a data-driven way from <u>unlabeled</u> <u>data</u> is one of the central problems in <u>unsupervised learning</u>.

AUTOENCODER

First of many simple but clever ideas: If we have inputs $\mathbf{x}_1, \ldots, \mathbf{x}_n \in \mathbb{R}^d$ but few or no targets y_1, \ldots, y_n , just make the inputs the targets.

- Let $f_{\boldsymbol{\theta}} : \mathbb{R}^d \to \mathbb{R}^d$ be our model.
- Let L_{θ} be a loss function. E.g. squared loss: $L_{\theta}(\mathbf{x}) = \left(\| \mathbf{x} - f_{\theta}(\mathbf{x}) \|_{2}^{2} \right) \longrightarrow \| \mathbf{x}_{i} - f_{\theta}(\mathbf{x}_{i}) \|_{2}^{2}$

• Train model:
$$\theta^* = \min_{\theta} \sum_{i=1}^n L_{\theta}(\mathbf{x})$$
.

If f_{θ} is a model that incorporates feature learning, then these features can be used for supervised tasks.

 f_{θ} is called an **autoencoder**. It maps input space to input space (e.g. images to images, french to french, PDE solutions to PDE solutions).

AUTOENCODER

Important property of autoencoders: no matter the architecture, there must always be a **bottleneck** with fewer parameters than the input. The bottleneck ensures information is "distilled" from low-level features to high-level features.



AUTOENCODER

Separately name the mapping from input to bottleneck and from bottleneck to output.

Encoder: $e : \mathbb{R}^d \to \mathbb{R}^k$



AUTOENCODER RECONSTRUCTION

Example image reconstructions from autoencoder:



https://www.biorxiv.org/content/10.1101/214247v1.full.pdf

Input parameters: (d = 49152)Bottleneck "latent" parameters: k = 1024. The best autoencoders do not work as well as supervised methods for feature extraction, but they require no labeled data.

There are a lot of cool applications of autoencoders beyond feature learning!

- Learned data compression.
- Denoising and in-painting.
- Data/image synthesis.

AUTOENCODERS FOR DATA COMPRESSION

Due to their bottleneck design, autoencoders perform **dimensionality reduction** and thus data compression.



Given input **x**, we can completely recover $f(\mathbf{x})$ from $\mathbf{z} = e(\mathbf{x})$. **z** typically has many fewer dimensions than **x** and for a typical

The best lossy compression algorithms are tailor made for specific types of data:

- JPEG 2000 for images
- MP3 for digital audio.
- MPEG-4 for video.

All of these algorithms take advantage of specific structure in these data sets. E.g. JPEG assumes images are locally "smooth".



engendining of addresses	C. Statement Products				
the providence of the second	0781307953	a an in the set of the	- POLINE CALE DE CALE	to repair many and	PERSONAL PROPERTY AND INCOME.
	1000100	na interio di serie	distriction in	NAME OF BRIDE	North Contern
		ant (b) a set (cash) () . B (data ca)			
and and a local designments			a many houses	a shares	and a state of the
	100000-000	earraine Scat			
Several and a province of	a a second	a shi ta mear a ta ar	- DO THE COUNTY	a	The second card
100000000000000000000000000000000000000		a lever and the second	REVERSORS	ISR CARDIN	NEW COLOR SCALE
LOOP DIVISION & CONTRACT, S. MINISTER, M.	CARLES MAN	CONTRACTOR CONTRACTOR	A PRIME PARAMAN		
e d'as l'alles as an i can ab	An infan	and and the second		a la fai han shada	e an e se par la se se se
THE R. LEWIS CO., Name of Address of the Party of the Par			THE REAL PROPERTY AND ADDRESS	in the second	
IN THE OWNER OF TAXABLE	02/0/20023	PROFILE STATE	m. Fri die their taken i		Interiore in Adver
STORE IC ADD. IN TAXABLE	Accounting	A RANK COLOR & SALAR	ALL	No. of ACCORDING STATE	CALCULATION DATE OF CALC
se alliniours in a					Contraction Contract
ALMAN TANILA OF AND MALIN	1011 (NJ 76	CAR INCOMEND & MARY	deline reacts provide a series	A WALL MANY NO. NA.	CURVERSE AND CONTRACT OF THE OWNER
and and balant sectors of spinsters	. A later	and group drawn		an when not	0.0131.000000
BICTORICS DOMESTIC DE LA DESERVICIÓN DE	SC SACOO CAIS	COLUMN TRANSPORT	C 111 23 0073 1894	NAME OF OCCUPANTS OF TAXABLE PARTY OF TA	
COLUMN TO BE AND ADDRESS OF	And Designed in the local division of the lo		CONTRACTOR OF		COLUMN TWO IS NOT THE OWNER.
and the late the second second	an out of the	CALCULATION OF THE OWNER OF THE O			a stricter of the strength

AUTOENCODERS FOR IMAGE COMPRESSION

With enough input data, autoencoders can be trained to find this structure on their own.



Proposed method, 5908 bytes (0.167 bit/px), PSNR: luma 23.38 dB/chroma 31.86 dB, MS-SSIM: 0.9219



JPEG 2000, 5908 bytes (0.167 bit/px), PSNR: luma 23.24 dB/chroma 31.04 dB, MS-SSIM: 0.8803



Proposed method, 6021 bytes (0.170 bit/px), PSNR: 24.12 dB, MS-SSIM: 0.9292



JPEG 2000, 6037 bytes (0.171 bit/px), PSNR: 23.47 dB, MS-SSIM: 0.9036

"End-to-end optimized image compression", Ballé, Laparra, Simoncelli

Need to be careful about how you choose loss function, design the network, etc. but can lead to much better image compression than "hand-tuned" algorithms like JPEG.

AUTOENCODERS FOR IMAGE CORRECTION



Image inpainting

Train autoencoder on <u>uncorrupted</u> images (unsupervised). Pass corrupted image **x** through autoencoder and return $f(\mathbf{x})$ as repaired result.

AUTOENCODERS LEARN COMPRESSED REPRESENTATIONS



Consider $128 \times 128 \times 3$ images with pixels values in 0, 1..., 255. How many possible images are there?

If **z** holds *k*, <u>8 bit</u> values, how many unique images **w** can be output by the autoencoder function *f*?

AUTOENCODERS LEARN COMPRESSED REPRESENTATIONS



An accurate autoencoder with a small bottleneck must have a representation space S that closely approximates I. Both will be much smaller than A.

AUTOENCODERS LEARN COMPRESSED REPRESENTATIONS



 $f(\mathbf{x}) = d(e(\mathbf{x}))$ projects an image **x** closer to the space of natural images.

Suppose we want to generate a random natural image. How might we do that?

• **Option 1**: Draw each pixel value in **x** uniformly at random. Draws a random image from *A*.



• **Option 2**: Draw **x** randomly from *S*, the space of images representable by the autoencoder.



How do we randomly select an image from \mathcal{S} ?



¹Some details to think about here. In reality, people use "variational autoencoders" (VAEs), which are a natural modification of AEs.

Deeper dive into understanding a simple, but powerful autoencoder architecture. Specifically we will view **principal component analysis (PCA)** as a type of autoencoder.

PCA is the "linear regression" of unsupervised learning: often the go-to baseline method for denoising, dimensionality reduction, etc.

Very important outside machine learning as well.



- One hidden layer. No non-linearity. No biases.
- Latent space of dimension k.
- Weight matrices are $\mathbf{W}_1 \in \mathbb{R}^{d \times k}$ and $\mathbf{W}_2 \in \mathbb{R}^{k \times d}$.

Given input $\mathbf{x} \in \mathbb{R}^d$, what is $f(\mathbf{x})$ expressed in linear algebraic terms? $(\mathbf{x} \leftarrow \mathbf{x})(\mathbf{z} \times \mathbf{k}) \rightarrow (\mathbf{z} \times \mathbf{k})$





Encoder:
$$z = \underline{e}(\underline{x}) = \underline{x}^T \underline{W}_1$$
. Decoder: $d(z) = \underline{z} \underline{W}_2$
 $\underline{x}^1 \omega, \ \omega_{\gamma}$

Given training data set $\underline{x_1}, \dots, \underline{x_n}$, let X denote our data matrix. Let $\tilde{X} = XW_1W_2$.



Natural squared autoencoder loss: Minimize $L(X, \tilde{X})$ where:

Goal: Find $\mathbf{W}_1, \mathbf{W}_2$ to minimize the Frobenius norm loss $\|\mathbf{X} - \tilde{\mathbf{X}}\|_F^2 = \|\mathbf{X} - \mathbf{X}\mathbf{W}_1\mathbf{W}_2\|_F^2$ (sum of squared entries).

LOW-RANK APPROXIMATION

Rank in linear algebra:

- The columns of a matrix with <u>column rank</u> k can all be written as linear combinations of just k columns.
- The rows of a matrix with <u>row rank k</u> can all be written as linear combinations of k rows. $\chi = (\chi \circ) v_{\mu}$



 \tilde{X} is a **low-rank matrix**. It only has rank k for $k \ll d$.

Principal component analysis is the task of finding W_1 , W_2 , which amounts to finding a rank k matrix \tilde{X} which approximates the data matrix X as closely as possible.

Finding the best **W**₁ and **W**₂ is a <u>non-convex</u> problem. We could try running an iterative method like gradient descent anyway. But there is also a direct algorithm!

SINGULAR VALUE DECOMPOSITION

Х



=



Where $\mathbf{U}^{\mathsf{T}}\mathbf{U} = \mathbf{I}, \mathbf{V}^{\mathsf{T}}\mathbf{V} = \mathbf{I}$, and $\underline{\sigma_1} \ge \underline{\sigma_2} \ge \dots \underline{\sigma_d} \ge 0$. I.e. **U** and **V** are orthogonal matrices.

This is called the **singular value decomposition**.

Can be computed in $O(nd^2)$ time (faster with approximation algos).

Let $\mathbf{u}_1, \ldots, \mathbf{u}_n \in \mathbb{R}^n$ denote the columns of **U**. I.e. the left singular vectors of **X**. Recall that orthogonality means that:



SINGULAR VALUE DECOMPOSITION



(Eckart-Young-Mirsky Theorem) For any $k \le d$, $\underline{X}_k = U_k \Sigma_k V_k^T$ is the optimal k rank approximation to X:

$$\mathbf{X}_k = \mathop{\mathrm{arg\,min}}_{\tilde{\mathbf{X}} ext{ with rank } \leq k} \| \underbrace{\mathbf{X} - \hat{\mathbf{X}}}_F \|_F^2.$$

That's great, but not quite in the form we wanted. Optimal rank k approximation is $\mathbf{X}_{k} = \mathbf{U}_{k} \mathbf{\Sigma}_{k} \mathbf{V}_{k}^{T}$. We want an approximation of the form:

$$\tilde{\mathbf{X}} = \mathbf{X} \mathbf{W}_1 \mathbf{W}_2$$



OPTIMAL LOW-RANK APPROXIMATION

Claim:
$$X_{k} = \bigcup_{k} \underbrace{XV_{k}}_{k} \underbrace{XV_{k}}_{k} \underbrace{XV_{k}}_{k} \underbrace{V_{k}}_{k} = x \underbrace{V_{k}}_{k} \underbrace{V_{k}} \underbrace{V_{k}}_{k} \underbrace{V_{k}}_{k} \underbrace{V_{$$



Usually **x**'s columns (features) are mean centered and normalized to variance 1 before computing principal components.

Computing the SVD.

- Full SVD: U,S,V = scipy.linalg.svd(X). Runs in O(nd²) time.
- Just the top k components:
 Uk, Sk, Vk = scipy.sparse.linalg.svds(X, k).
 Runs in roughly O(ndk) time.

Recall that for a matrix $\underline{\mathbf{M}} \in \mathbb{R}^{p \times p}$, $\underline{\mathbf{q}}$ is an <u>eigenvector</u> of \mathbf{M} if $\lambda \underline{\mathbf{q}} = \underline{\mathbf{M}} \underline{\mathbf{q}}$ for any scalar λ .

- (U) <u>columns</u> (the left singular vectors) are the orthonormal eigenvectors of (XX^T)
- V_{S} columns (the right singular vectors) are the orthonormal eigenvectors of $X^T X$.

$$\cdot \sigma_i^2 = \lambda_i (\mathbf{X}\mathbf{X}^T) = \lambda_i (\mathbf{X}^T\mathbf{X})$$

Exercise: Verify this directly. This means you can use any eigensolver for computing the SVD.

PCA APPLICATIONS

Like any autoencoder, PCA can be used for:

- Feature extraction
- Denoising and rectification
- Data generation
- Compression
- Visualization




LOW-RANK APPROXIMATION



36

Error vs. k is dictated by X's singular values. The singular values are often called the **spectrum** of X. $6\sqrt{7}$

$$\|\mathbf{X} - \mathbf{X}_k\|_F^2 = \sum_{i=k+1}^d \sigma_i^2.$$



»6u

COLUMN REDUNDANCY

Colinearity of data features leads to an approximately low-rank data matrix.



sale price $\approx 1.05 \cdot \text{list}$ price. property tax $\approx .01 \cdot \text{list}$ price. Sometimes these relationships are simple, other times more complex. But as long as there exists <u>linear</u> relationships between features, we will have a lower rank matrix.

yard size
$$\approx$$
 lot size $-\frac{1}{2}$ · square footage.

cumulative GPA
$$\approx \frac{1}{4} \cdot$$
 year 1 GPA $+ \frac{1}{4} \cdot$ year 2 GPA
 $+ \frac{1}{4} \cdot$ year 3 GPA $+ \frac{1}{4} \cdot$ year 4 GPA.

LOW-RANK INTUITION

Two other examples of data with good low-rank approximations:

- 1. Genetic data: single nucleotide polymorphisms (SNPs) loci 144 312 436 800 943 indivi<u>dual</u>1 A T T T G G G individual 2 C A individual n Т Α G
- 2. "Term-document" matrix with bag-of-words data:

	c	ar "	oan 'C	N _{SO}		•••		0	700	Cat
	doc_1	0	0	1	0	0	1	1	0	0
	doc_2	0	0	0	1	0	1	0	0	0
	:	1	1	0	1	0	0	0	1	0
l=EP	•	0	0	0	0	0	0	0	1	1
	doc_n	1	0	0	0	0	0	0	1	1

PRINCIPAL COMPONENTS

3:5 Contra

What do principal components and loading vectors look like?

PRINCIPAL COMPONENTS

MNIST principal components:



Often principal components are difficult to interpret.

What do the loading vectors looks like?

The loading vector **z** for an example **x** contains coefficients which recombine the top *k* principal components v_1, \ldots, v_k to approximately reconstruct **x**.



Provide a short <u>"finger prin</u>t" for any image **x** which can be used to reconstruct that image.

For any **x** with loading vector **z**, the *i*th entry z_i is the inner product similarity between **x** and the *i*th principal component, \mathbf{v}_i .



LOADING VECTORS: PROJECTION VIEW



Since $\mathbf{v}_1, \ldots, \mathbf{v}_k$ are orthonormal, this operation is a projection onto first *k* principal components.

I.e. we are projecting **x** onto the *k*-dimensional subspace spanned by $\mathbf{v}_1, \ldots, \mathbf{v}_k$.

For an example \mathbf{x}_i , the loading vector \mathbf{z}_i contains the coordinates in the projection space:



Important takeaway for data visualization and more: Latent feature vectors preserve similarity and distance information in the original data.

Let $\mathbf{x}_1 \dots, \mathbf{x}_n \in \mathbb{R}^d$ be our original data vectors, $\mathbf{z}_1 \dots, \mathbf{z}_n \in \mathbb{R}^k$ be our loading vectors (encoding), and $\mathbf{\tilde{x}}_1 \dots, \mathbf{\tilde{x}}_n \in \mathbb{R}^d$ be our low-rank approximated data.

We have:

$$\begin{split} \|\tilde{\mathbf{x}}_i\|_2^2 &= \|\mathbf{z}_i\|_2^2\\ \langle \underline{\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j} \rangle &= \langle \mathbf{z}_i, \mathbf{z}_j \rangle\\ \|\underline{\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j}\|_2^2 &= \|\mathbf{z}_i - \mathbf{z}_j\|_2^2 \end{split}$$

Conclusion: If our data had a good low rank approximation, i.e. $\|\tilde{\mathbf{x}}_{i}\|_{2}^{2} \approx \|\mathbf{x}_{i}\|_{2}^{2} \langle \tilde{\mathbf{x}}_{i}, \tilde{\mathbf{x}}_{j} \rangle \approx \langle \mathbf{x}_{i}, \mathbf{x}_{j} \rangle$ and $\|\tilde{\mathbf{x}}_{i} - \tilde{\mathbf{x}}_{j}\|_{2}^{2} \approx \|\mathbf{x}_{i} - \mathbf{x}_{j}\|_{2}^{2}$, we expect that: $\sum_{i=1}^{n} \langle \mathcal{V}_{i}, \mathcal{V}_{i} \rangle \langle \mathbf{x}_{i}, \mathbf{x}_{j} \rangle \approx \langle \mathbf{z}_{i}, \mathbf{z}_{j} \rangle$ $\|\mathbf{x}_{i}\|_{2}^{2} \approx \|\mathbf{z}_{i}\|_{2}^{2}$ $\|\mathbf{x}_{i} - \mathbf{x}_{j}\|_{2}^{2} \approx \|\mathbf{z}_{i}\|_{2}^{2}$

Useful in obtaining short "finger prints" for complex data.

Note: this is not true of most autoencoders, but unique to PCA. Typically compressions themselves cannot be directly used to approximate distance, similarity,. etc.

TERM DOCUMENT MATRIX



Word-document matrices tend to be low rank.

Documents tend to fall into a relatively small number of different categories, which use similar sets of words:

- Financial news: markets, analysts, dow, rates, stocks
- US Politics: president, senate, pass, slams, twitter, media
- StackOverflow posts: python, help, convert, javascript

LATENT SEMANTIC ANALYSIS

Latent semantic analysis = PCA applied to a word-document matrix (usually from a large corpus). One of the most fundamental techniques in **natural language processing** (NLP).



Each column of **z** corresponds to a latent "category" or "topic". Corresponding row in **Y** corresponds to the "frequency" with which different words appear in documents on that topic.

LATENT SEMANTIC ANALYSIS

(21,2) ~ (x1, Kj)

Similar documents have similar LSA document vectors. I.e. $\langle \boldsymbol{z}_i, \boldsymbol{z}_j \rangle$ is large.

- z_i provides a more compact "finger print" for documents than the long bag-of-words vectors. Useful for e.g search engines.
- Comparing document vectors is often <u>more effective</u> than comparing raw BOW features. Two documents can have ((z_i, z_j))large even if they have no overlap in words. E.g. because both share a lot of words with words with another document k, or with a bunch of other documents.

Same fingerprinting idea was also important in early facial recognition systems based on "eigenfaces":



Each image above is one of the principal components of a dataset containing images of faces.

SEMANTIC EMBEDDINGS

Document embeddings are clearly useful. What about the word embeddings? It turns out these are super useful as well!

Reminder: The *i*, *j* entry of \tilde{X} equals $\langle z_i, y_j \rangle$.





If two words often appear in the same documents, their word vectors tend to point more in the same direction.

Result: Map words to numerical vectors in a <u>semantically</u> meaningful way. Similar words map to similar vectors. Dissimilar words to dissimilar vectors.



Extremely useful "side-effect" of LSA.

Capture e.g. the fact that "great" and "excellent" are near synonyms. Or that "difficult" and "easy" are antonyms.

For similar words, $\langle \mathbf{y}_i, \mathbf{y}_j \rangle$ should be large. I.e. \mathbf{y}_i and \mathbf{y}_j point in the same direction.



'**Review 1:** Very small and handy for traveling or camping. Excellent quality, operation, and appearance.

Review 2: So far this thing is great. Well designed, compact, and easy to use. I'll never use another can opener.

Review 3: Not entirely sure this was worth \$20. Mom couldn't figure out how to use it and it's fairly difficult to turn for someone with arthritis.

Goal is to classify reviews as "positive" or "negative".

Vocabulary: Small, handy, excellent, great, quality, compact, easy, difficult.

Review 1: Very small and handy for traveling or camping. Excellent quality, operation, and appearance.

[, , , , , , , , ,]

Review 2: So far this thing is great. Well designed, compact, and easy to use. I'll never use another can opener.

[, , , , , , , ,]

Review 3: Not entirely sure this was worth \$20. Mom couldn't figure out how to use it and it's fairly difficult to turn for someone with arthritis.

[, , , , , , , , , , , , , ,]

Bag-of-words approach typically only works for large data sets.

The features do not capture the fact that "great" and "excellent" are near synonyms. Or that "difficult" and "easy" are antonyms.



This can be addressed by first mapping words to <u>semantically</u> <u>meaningful vectors</u>. That mapping can be trained using a much large corpus of text than the data set you are working with (e.g. Wikipedia, Twitter, news data sets). How to go from word embeddings to features for a whole sentence or chunk of text?



USING WORD EMBEDDINGS



To avoid issues with inconsistent sentence length, word ordering, etc., can concatenate a fixed number of top <u>principal</u> <u>components</u> of the matrix of word vectors:



There are much more complicated approaches that account for word position in a sentence. Lots of pretrained libraries available (e.g. Facebook's **\InferSent**).

62



We chose Z to equal $\underline{XV}_k = U_k \Sigma_k$ and $\mathbf{Y} = \underline{V}_k^T$. Could have just as easily set $\mathbf{Z} = U_k$ and $\mathbf{Y} = \boldsymbol{\Sigma}_k \mathbf{V}_k^T$, so Z has orthonormal columns.





What does the i, j entry of $X^T X$ reprent?

 $\langle \mathbf{y}_i, \mathbf{y}_j \rangle$ is <u>larger</u> if *word*_i and *word*_j appear in more documents together (high value in **word-word co-occurrence matrix**, $\mathbf{X}^T \mathbf{X}$). Similarity of word embeddings mirrors similarity of word context.

General word embedding recipe:

- 1. (Choose similarity metric $k(word_i, word_j)$ which can be computed for any pair of words.
- 2. Construct similarity matrix $\mathbf{M} \in \mathbb{R}^{n \times n}$ with $\mathbf{M}_{i,j} = k(\underline{word}_i, word_j)$.

3. Find low rank approximation $\underline{M} \approx \underline{Y}^T \underline{Y}$ where $\underline{Y} \in \mathbb{R}^{k \times n}$.

4. Columns of $\widehat{(Y)}$ are word embedding vectors.

We expect that $\langle \underline{y}_i, \underline{y}_j \rangle$ will be larger for more similar words.



How do current state-of-the-art methods differ from LSA?

- Similarity based on co-occurrence in smaller chunks of words. E.g. in sentences or in any consecutive sequences of 3, 4, or 10 words.
- Usually transformed in non-linear way. E.g. $k(word_i, word_j) = \frac{(p(i,j))}{p(i)p(j)}$ where p(i,j) is the frequency both i, j appeared together, and p(i), p(j) is the frequency either one appeared.

Computing word similarities for "window size" 4:

The girl walks to her dog to the park. It can take a long time to park your car in NYC. The dog park is always crowded on Saturdays.

The girl walks to her dog to the park. It can take a long time to park your car in NYC. The dog park is always crowded on Saturdays.

The girl walks to her dog to the park. It can take a long time to park your car in NYC. The dog park is always crowded on Saturdays.

	dog	park	park crowded		
gop	0	2	0	3	
park	2	0	1	2	
crowded	0	1	0	0	
the	3	2	0	0	

Current state of the art models: GloVE, word2vec.

- word2vec was originally presented as a shallow neural network model, but it is equivalent to matrix factorization method (Levy, Goldberg 2014).
- For word2vec, similarity metric is the "point-wise mutual information": $\log \frac{p(i,j)}{p(i)p(j)}$.

CAVEAT ABOUT FACTORIZATION



SVD will not return a symmetric factorization in general. In fact, if **M** is not positive semidefinite² then the optimal low-rank approximation does not have this form.

²I.e., *k*(*word*_{*i*}, *word*_{*i*}) is not a positive semidefinite kernel.
CAVEAT ABOUT FACTORIZATION



- For each word *i* we get a left and right embedding vector
 w_i and y_i. It's reasonable to just use one or the other.
- If (y_i, y_j) is large and positive, we expect that y_i and y_j have similar similarity scores with other words, so they typically are still related words.
- Another option is to use as your features for a word the concatenation [w_i, y_i]

Lots of <u>pre-trained</u> word vectors are available online:

- Original gloVe website: https://nlp.stanford.edu/projects/glove/.
- Compilation of many sources: https://github.com/3Top/word2vec-api

Lots of cool demos for what can be done with these embeddings. E.g. "vector math" to solve analogies.



FORWARD LOOKING APPLICATION: UNSUPERVISED TRANSLATION



- Train word-embeddings for languages separately. Obtain lowish dimensional point clouds of words.
- Perform rotation/alignment to match up these point clouds.
- Equivalent words should land on top of each other.

No needs for labeled training data like translated pairs of sentences!

Why not monkey or whale language?



Earth Species Project (www.earthspecies.org), CETI Project (www.projectceti.org)

The same approach used for word embeddings can be used to obtain meaningful numerical features for any other data where there is a natural notion of similarity.



For example, the items could be nodes in a social network graph. Maybe be want to predict an individuals age, level of interest in a particular topic, political leaning, etc.

NODE EMBEDDINGS



Generate random walks (e.g. "sentences" of nodes) and measure similarity by node co-occurence frequency.



NODE EMBEDDINGS

Again typically normalized and apply a non-linearity (e.g. log) as in word embeddings.



Popular implementations: **DeepWalk**, **Node2Vec**. Again initially derived as simple neural network models, but are equivalent to matrix-factorization (Qiu et al. 2018).

BIMODAL EMBEDDINGS

We can also create embeddings that represent different types of data. OpenAl's clip architecture:



Goal: Train embedding architectures so that $\langle T_i, I_j \rangle$ are similar if image and sentence are similar.

CLIP TRAINING

What do we use as ground truth similarities during training? Sample a batch of sentence/image pairs and just use identity matrix.



My new puppy!	1	0	0
Best dim sum ever.	0	1	0
NYC in the rain.	0	0	1

This is called <u>contrastive learning</u>. Train unmatched text/image pairs to have nearly orthogonal embedding vectors.

CLIP FOR ZERO-SHOT LEARNING

Learning Transferable Visual Models From Natural Language Supervision

Alec Radford^{*1} Jong Wook Kim^{*1} Chris Hallacy¹ Aditya Ramesh¹ Gabriel Goh¹ Sandhini Agarwal¹ Girish Sastry¹ Amanda Askell¹ Pamela Mishkin¹ Jack Clark¹ Gretchen Krueger¹ Ilya Sutskever¹



IMAGE SYNTHESIS

AUTOENCODERS LEARN COMPRESSED REPRESENTATIONS



 $f(\mathbf{x}) = d(e(\mathbf{x}))$ projects an image **x** closer to the space of natural images.

Suppose we want to generate a random natural image. How might we do that?

• **Option 1**: Draw each pixel value in **x** uniformly at random. Draws a random image from *A*.



• **Option 2**: Draw **x** randomly from *S*, the space of images representable by the autoencoder.



How do we randomly select an image from \mathcal{S} ?

Autoencoder approach to generative ML: Feed random inputs into decode to produce random realistic outputs.





AUTOENCODERS FOR DATA GENERATION



Variational auto-encoders attempt to resolve this issue.

Variational auto-encoders attempt to resolve this issue. Basic ideas:

- Add noise during training.
- Add penalty term so that distribution of code vectors generated looks like mean 0, variance 1 Gaussian.



Variation AE's give very good results, but tends to produce images with immediately recognizable flaws (e.g. soft edges, high-frequency artifacts).





Lots of efforts to hand-design regularizers that penalize images that don't look realisitic to the human eye. **Main idea behind GANs:** Use machine learning to automatically encourage realistic looking images.

 $\min_{\theta} L(\theta) + P(\theta)$

GENERATIVE ADVERSARIAL NETWORKS (GANS)



Let x_1, \ldots, x_n be real images and let z_1, \ldots, z_m be random code vectors. The goal of the discriminator is to output a number between [0, 1] which is close to 0 if the image is fake, close to 1 if it's real.

Train weights of discriminator D_{θ} to minimize:

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^{n} -\log\left(D_{\boldsymbol{\theta}}(\mathbf{x}_{i})\right) + \sum_{i=1}^{m} -\log\left(1 - D_{\boldsymbol{\theta}}(G_{\boldsymbol{\theta}'}(\mathbf{z}_{i}))\right)$$

GENERATIVE ADVERSARIAL NETWORKS (GANS)



Goal of the generator $G_{\theta'}$ is the opposite. We want to maximize: $\max_{\theta'} \sum_{i=1}^{} -\log(1 - D_{\theta}(G_{\theta'}(\mathbf{z}_i)))$

This is called an "adversarial loss function". *D* is playing the role of the adversary.

$$\boldsymbol{\theta}^*, \boldsymbol{\theta}'^* \text{ solve } \min_{\boldsymbol{\theta}} \max_{\boldsymbol{\theta}'} \sum_{i=1}^n -\log \left(D_{\boldsymbol{\theta}}(\mathbf{x}_i) \right) + \sum_{i=1}^m -\log \left(1 - D_{\boldsymbol{\theta}}(G_{\boldsymbol{\theta}'}(\mathbf{z}_i)) \right)$$

This is called a minimax optimization problem. <u>Really tricky to</u> solve in practice.

- Repeatedly play: Fix one of θ* or θ'*, train the other to convergence, repeat.
- Simultaneous gradient descent: Run a single gradient descent step for each of θ^{*}, θ'^{*} and update D and G accordingly. Difficult to balance learning rates.
- Lots of tricks (e.g., slightly different loss functions) can help.

State of the art until a few years ago.



DIFFUSION

Auto-encoder/GAN approach: Input noise, map directly to image.

Diffusion: Slowly move from noise to image.



We will post a demo for generating digits by training on MNIST.





SEMANTIC EMBEDDINGS + DIFFUSION

Text to image synthetsis: Dall-E, Imagen, Stable Diffusion



"A chair that looks like an avocado"