

Lecture 7: Representation Learning

André Martins, Francisco Melo, Mário Figueiredo



Deep Learning Course, Fall 2021

Announcements

Deadline for Homework #1 is this **Wednesday end of day!**

- Please submit your solutions and code in Fenix.
- No late days allowed!!
- Solutions will be posted the day after.

Homework #2 will be posted this Wednesday!

- Deadline Jan 31.
- Start early!!

Today's Roadmap

Today's lecture is about:

- Representation learning.
- Principal component analysis (PCA) and auto-encoders.
- Denoising auto-encoders.
- Distributed representations.
- Word embeddings and negative sampling.
- Multilingual and contextual word embeddings.

Outline

① Representation Learning

Hierarchical Compositionality

Distributed Representations

Auto-Encoders

Word Embeddings

② Conclusions

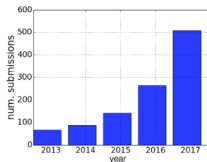
Representations

- A key feature of NNs is their ability to **learn representations** $\phi(x)$
- Standard linear models require **manually engineered features** $\phi(x)$
- **Representations** are useful for several reasons:
 - (i) They can make our models **more expressive and more accurate**
 - (ii) They may allow **transferring** representations from one task to another
- We talked about (i) when discussing the multi-layer perceptron
- In this lecture, we'll focus on (ii)

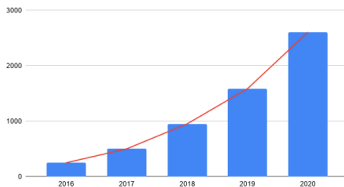
Representation Learning

This is becoming an extremely popular topic!

Number of submissions to the “International Conference on Learning Representations” (ICLR):



ICLR Submissions



Representation learning almost became a synonym of deep learning

Outline

① Representation Learning

Hierarchical Compositionality

Distributed Representations

Auto-Encoders

Word Embeddings

② Conclusions

Hierarchical Compositionality

Key Idea: deep(er) NNs learn **coarse-to-fine** representation layers.

Vision:

- pixels → edges → textures → motifs → parts → objects → scenes

Speech:

- audio samples → spectral bands → formants → motifs → phonemes → words

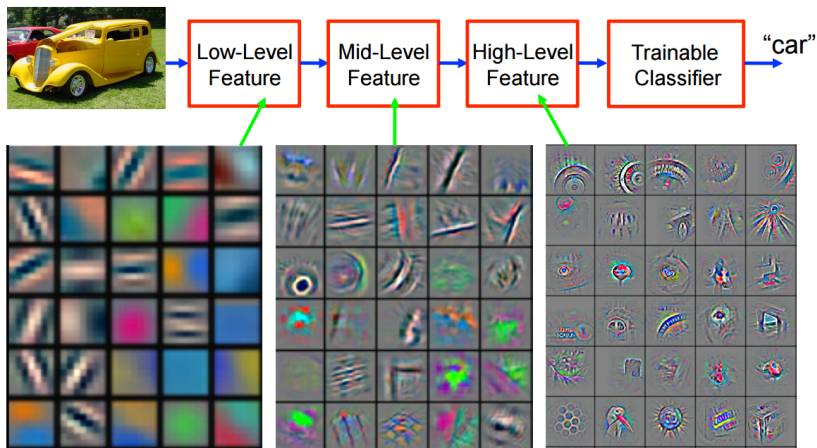
Text:

- characters → words → phrases → sentences → stories

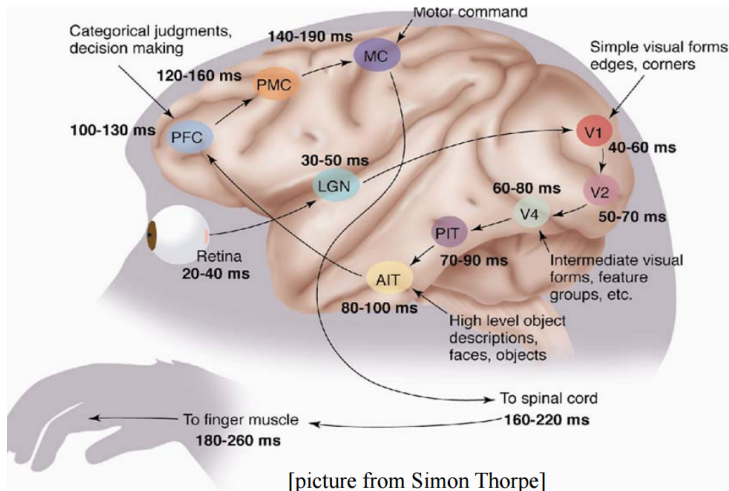
(Slide inspired by Marc'Aurelio Ranzato and Yann LeCun)

Hierarchical Compositionality

Feature visualization of convolutional NNs trained on ImageNet (Zeiler and Fergus, 2013):



The Mammalian Visual Cortex is Hierarchical



(Slide inspired by Marc'Aurelio Ranzato and Yann LeCun)

What's in Each Layer

- Bottom level layers (closer to inputs) tend to learn **low-level representations** (corners, edges)
- Upper level layers (farther away from inputs) learn **more abstract representations** (shapes, forms, objects)

This holds for images, text, etc.

Outline

① Representation Learning

Hierarchical Compositionality

Distributed Representations

Auto-Encoders

Word Embeddings

② Conclusions

Distributed Representations (Hinton, 1984)

This is a central concept in neural networks.

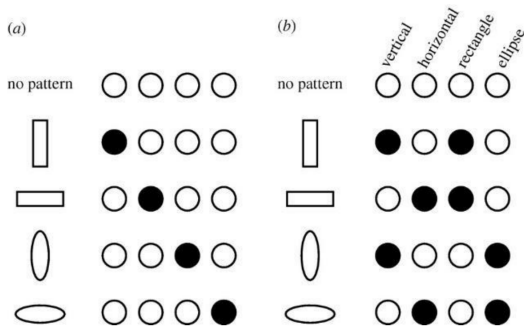
Key questions:

- How can a NN so effectively **represent objects**, if it has only a few **hidden units** (i.e. much fewer than possible objects)?
- What is each **hidden unit** actually **representing**?
- How can a NN **generalize** to objects that is has not seen before?

Local vs Distributed Representations

Consider two alternative representations:

- **Local** (one-hot) representations (one dimension per object)
- **Distributed representations** (one dimension per **property**)



(Slide inspired by Moontae Lee and Dhruv Batra)

Distributed Representations

Key idea: no single neuron “encodes” everything; groups of neurons (e.g. in the same hidden layer) work together!

Distributed Representations

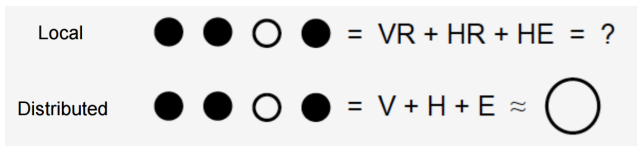
Key idea: no single neuron “encodes” everything; groups of neurons (e.g. in the same hidden layer) work together!

cf. the [grandmother cell](#)



The Power of Distributed Representations

- Distributed representations are **more compact** (there can be $O(\exp N)$ objects combining N properties)
- They are also **more powerful**, as they can generalize to unseen objects in a meaningful way:



(Slide inspired by Moontae Lee and Dhruv Batra)

The Power of Distributed Representations

- For this to work, hidden units should capture **diverse properties** of objects (not all capturing the same property)
- Usually ensured by random initialization of the weights
- Initializing all the units to the same weights, we would never break the symmetry!
- **Side note:** a NN computes the same function if we permute the hidden units within the same layer (order doesn't matter, only diversity)

Next: how to learn useful object representations from raw inputs (no labels)?

Example: Unsupervised Pre-Training

- Training deep NNs (with many hidden layers) can be challenging
- This has been a major difficulty with NNs for a long time
- Initialize hidden layers using **unsupervised learning** (Erhan et al., 2010):
 - Force network to **represent latent structure** of input distribution
 - Encourage hidden layers to **encode** that structure
 - This can be done with an **auto-encoder**!

Outline

① Representation Learning

Hierarchical Compositionality

Distributed Representations

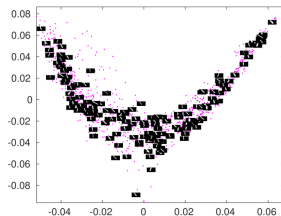
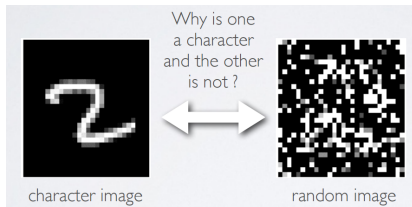
Auto-Encoders

Word Embeddings

② Conclusions

Data Manifold

Key idea: learn the **manifold** where the input objects live



(Image credit: Hugo Larochelle)

Learn representations that encode well points in that **manifold**

Auto-Encoders

Auto-encoder: feed-forward NN trained to reproduce its input at the output

Encoder:

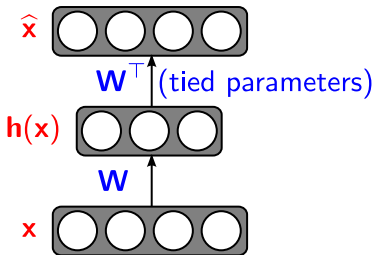
$$h(x) = g(Wx + b)$$

Decoder:

$$\hat{x} = W^T h(x) + c$$

Loss function (for real-valued inputs):

$$L(\hat{x}; x) = \frac{1}{2} \|\hat{x} - x\|^2$$



The Simplest Auto-Encoder

What happens if the activation function g is linear?

The Simplest Auto-Encoder

What happens if the activation function g is linear?

Principal Component Analysis (PCA)!

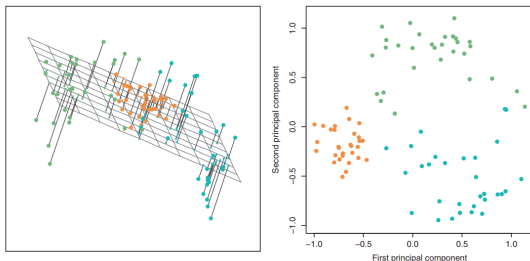


FIGURE 10.2. *Ninety observations simulated in three dimensions. Left: the first two principal component directions span the plane that best fits the data. It minimizes the sum of squared distances from each point to the plane. Right: the first two principal component score vectors give the coordinates of the projection of the 90 observations onto the plane. The variance in the plane is maximized.*

(From “An Introduction to Statistical Learning” by James, Witten, Hastie, Tibshirani)

Interlude: PCA

Consider N points in \mathbb{R}^D : $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$

Interlude: PCA

Consider N points in \mathbb{R}^D : $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$

Goal: find **good** K -dimensional representation of the points, with $K < D$,

$$\mathbf{x}^{(i)} \simeq \hat{\mathbf{x}}^{(i)} = \sum_{j=1}^K \alpha_j^{(i)} \mathbf{u}_j,$$

where $\{\mathbf{u}_1, \dots, \mathbf{u}_K\} \subset \mathbb{R}^D$ is an **orthonormal basis** of a **subspace** of \mathbb{R}^D .

Interlude: PCA

Consider N points in \mathbb{R}^D : $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$

Goal: find **good** K -dimensional representation of the points, with $K < D$,

$$\mathbf{x}^{(i)} \simeq \hat{\mathbf{x}}^{(i)} = \sum_{j=1}^K \alpha_j^{(i)} \mathbf{u}_j,$$

where $\{\mathbf{u}_1, \dots, \mathbf{u}_K\} \subset \mathbb{R}^D$ is an **orthonormal basis** of a **subspace** of \mathbb{R}^D .

If $\{\mathbf{u}_1, \dots, \mathbf{u}_K\}$ is fixed, and the approximation is in Euclidean norm,

$$\hat{\mathbf{x}}^{(i)} = \arg \min_{\mathbf{x}'} \|\mathbf{x}^{(i)} - \mathbf{x}'\|_2^2, \quad \text{subject to } \mathbf{x}' \in \text{span}(\{\mathbf{u}_1, \dots, \mathbf{u}_K\})$$

Interlude: PCA

Consider N points in \mathbb{R}^D : $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$

Goal: find **good** K -dimensional representation of the points, with $K < D$,

$$\mathbf{x}^{(i)} \simeq \hat{\mathbf{x}}^{(i)} = \sum_{j=1}^K \alpha_j^{(i)} \mathbf{u}_j,$$

where $\{\mathbf{u}_1, \dots, \mathbf{u}_K\} \subset \mathbb{R}^D$ is an **orthonormal basis** of a **subspace** of \mathbb{R}^D .

If $\{\mathbf{u}_1, \dots, \mathbf{u}_K\}$ is fixed, and the approximation is in Euclidean norm,

$$\hat{\mathbf{x}}^{(i)} = \arg \min_{\mathbf{x}'} \|\mathbf{x}^{(i)} - \mathbf{x}'\|_2^2, \quad \text{subject to } \mathbf{x}' \in \text{span}(\{\mathbf{u}_1, \dots, \mathbf{u}_K\})$$

has a well-known solution: **orthogonal projection**,

$$\hat{\mathbf{x}}^{(i)} = \sum_{j=1}^K \mathbf{u}_j \mathbf{u}_j^T \mathbf{x}^{(i)} = \mathbf{U} \mathbf{U}^T \mathbf{x}^{(i)}$$

where $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_K] \in \mathbb{R}^{D \times K}$; notice that $\mathbf{U}^T \mathbf{U} = \mathbf{I}$.

Interlude: PCA

Minimizing the average error for the N points in \mathbb{R}^D : $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$:

$$\min_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}^{(i)} - \mathbf{U} \mathbf{U}^T \mathbf{x}^{(i)}\|_2^2$$

Interlude: PCA

Minimizing the average error for the N points in \mathbb{R}^D : $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$:

$$\min_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}^{(i)} - \mathbf{U} \mathbf{U}^T \mathbf{x}^{(i)}\|_2^2$$

Notice that

$$\|\mathbf{x}^{(i)} - \mathbf{U} \mathbf{U}^T \mathbf{x}^{(i)}\|_2^2 = ((\mathbf{x}^{(i)})^T - (\mathbf{x}^{(i)})^T \mathbf{U} \mathbf{U}^T) (\mathbf{x}^{(i)} - \mathbf{U} \mathbf{U}^T \mathbf{x}^{(i)})$$

$$\stackrel{(a)}{=} \|\mathbf{x}^{(i)}\|_2^2 - (\mathbf{x}^{(i)})^T \mathbf{U} \mathbf{U}^T \mathbf{x}^{(i)}$$

$$= \|\mathbf{x}^{(i)}\|_2^2 - \sum_{j=1}^K \underbrace{(\mathbf{x}^{(i)})^T \mathbf{u}_j}_{\text{scalar}} \underbrace{\mathbf{u}_j^T \mathbf{x}^{(i)}}_{\text{scalar}}$$

indep. of \mathbf{U}

$$= \overbrace{\|\mathbf{x}^{(i)}\|_2^2}^{\text{indep. of } \mathbf{U}} - \sum_{j=1}^K \mathbf{u}_j^T \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T \mathbf{u}_j$$

$$(a) \text{ since } -2(\mathbf{x}^{(i)})^T \mathbf{U} \mathbf{U}^T \mathbf{x}^{(i)} + \underbrace{(\mathbf{x}^{(i)})^T \mathbf{U} \mathbf{U}^T \mathbf{U} \mathbf{U}^T \mathbf{x}^{(i)}}_{\mathbf{I}} = -(\mathbf{x}^{(i)})^T \mathbf{U} \mathbf{U}^T \mathbf{x}^{(i)}$$

Interlude: PCA

Minimizing the average error for the N points in \mathbb{R}^D : $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$:

$$\arg \min_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}^{(i)} - \mathbf{U} \mathbf{U}^T \mathbf{x}^{(i)}\|_2^2 = \arg \max_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K \mathbf{u}_j^T \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T \mathbf{u}_j$$

Interlude: PCA

Minimizing the average error for the N points in \mathbb{R}^D : $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$:

$$\begin{aligned} \arg \min_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}^{(i)} - \mathbf{U} \mathbf{U}^T \mathbf{x}^{(i)}\|_2^2 &= \arg \max_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K \mathbf{u}_j^T \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T \mathbf{u}_j \\ &= \arg \max_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \sum_{j=1}^K \mathbf{u}_j^T \underbrace{\left(\frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T \right)}_{\hat{\Sigma}} \mathbf{u}_j \end{aligned}$$

Interlude: PCA

Minimizing the average error for the N points in \mathbb{R}^D : $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$:

$$\begin{aligned}\arg \min_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}^{(i)} - \mathbf{U} \mathbf{U}^T \mathbf{x}^{(i)}\|_2^2 &= \arg \max_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K \mathbf{u}_j^T \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T \mathbf{u}_j \\ &= \arg \max_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \sum_{j=1}^K \mathbf{u}_j^T \underbrace{\left(\frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T \right)}_{\hat{\Sigma}} \mathbf{u}_j \\ &= \arg \max_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \sum_{j=1}^K \mathbf{u}_j^T \hat{\Sigma} \mathbf{u}_j\end{aligned}$$

$\hat{\Sigma}$ is the sample covariance, assuming **centred** data: $\sum_{i=1}^N \mathbf{x}^{(i)} = 0$.

PCA

Let's start with $K = 1$,

$$\mathbf{u}_1 = \arg \max_{\|\mathbf{u}\|_2^2=1} \mathbf{u}^T \hat{\Sigma} \mathbf{u}$$

PCA

Let's start with $K = 1$,

$$\mathbf{u}_1 = \arg \max_{\|\mathbf{u}\|_2^2=1} \mathbf{u}^T \hat{\Sigma} \mathbf{u}$$

Lagrangian:

$$L(\mathbf{u}, \lambda) = \mathbf{u}^T \hat{\Sigma} \mathbf{u} + \lambda(1 - \|\mathbf{u}\|_2^2)$$

PCA

Let's start with $K = 1$,

$$\mathbf{u}_1 = \arg \max_{\|\mathbf{u}\|_2^2=1} \mathbf{u}^T \hat{\Sigma} \mathbf{u}$$

Lagrangian:

$$L(\mathbf{u}, \lambda) = \mathbf{u}^T \hat{\Sigma} \mathbf{u} + \lambda(1 - \|\mathbf{u}\|_2^2)$$

Setting the gradient to zero: $2\hat{\Sigma}\mathbf{u} - 2\lambda\mathbf{u} = 0 \Rightarrow \hat{\Sigma}\mathbf{u} = \lambda\mathbf{u}$,

PCA

Let's start with $K = 1$,

$$\mathbf{u}_1 = \arg \max_{\|\mathbf{u}\|_2^2=1} \mathbf{u}^T \hat{\Sigma} \mathbf{u}$$

Lagrangian:

$$L(\mathbf{u}, \lambda) = \mathbf{u}^T \hat{\Sigma} \mathbf{u} + \lambda(1 - \|\mathbf{u}\|_2^2)$$

Setting the gradient to zero: $2\hat{\Sigma}\mathbf{u} - 2\lambda\mathbf{u} = 0 \Rightarrow \hat{\Sigma}\mathbf{u} = \lambda\mathbf{u}$,

i.e., \mathbf{u} is an eigenvector of $\hat{\Sigma}$. Plugging above:

$$\mathbf{u}^T \hat{\Sigma} \mathbf{u} = \mathbf{u}^T \mathbf{u} \lambda = \|\mathbf{u}\|_2^2 \lambda = \lambda$$

PCA

Let's start with $K = 1$,

$$\mathbf{u}_1 = \arg \max_{\|\mathbf{u}\|_2^2=1} \mathbf{u}^T \hat{\Sigma} \mathbf{u}$$

Lagrangian:

$$L(\mathbf{u}, \lambda) = \mathbf{u}^T \hat{\Sigma} \mathbf{u} + \lambda(1 - \|\mathbf{u}\|_2^2)$$

Setting the gradient to zero: $2\hat{\Sigma}\mathbf{u} - 2\lambda\mathbf{u} = 0 \Rightarrow \hat{\Sigma}\mathbf{u} = \lambda\mathbf{u}$,

i.e., \mathbf{u} is an eigenvector of $\hat{\Sigma}$. Plugging above:

$$\mathbf{u}^T \hat{\Sigma} \mathbf{u} = \mathbf{u}^T \mathbf{u} \lambda = \|\mathbf{u}\|_2^2 \lambda = \lambda$$

Conclusion: \mathbf{u}_1 is the eigenvector of the largest eigenvalue of $\hat{\Sigma}$.

PCA

Let's start with $K = 1$,

$$\mathbf{u}_1 = \arg \max_{\|\mathbf{u}\|_2^2=1} \mathbf{u}^T \hat{\Sigma} \mathbf{u}$$

Lagrangian:

$$L(\mathbf{u}, \lambda) = \mathbf{u}^T \hat{\Sigma} \mathbf{u} + \lambda(1 - \|\mathbf{u}\|_2^2)$$

Setting the gradient to zero: $2\hat{\Sigma}\mathbf{u} - 2\lambda\mathbf{u} = 0 \Rightarrow \hat{\Sigma}\mathbf{u} = \lambda\mathbf{u}$,

i.e., \mathbf{u} is an eigenvector of $\hat{\Sigma}$. Plugging above:

$$\mathbf{u}^T \hat{\Sigma} \mathbf{u} = \mathbf{u}^T \mathbf{u} \lambda = \|\mathbf{u}\|_2^2 \lambda = \lambda$$

Conclusion: \mathbf{u}_1 is the eigenvector of the largest eigenvalue of $\hat{\Sigma}$.

Notice: $\hat{\Sigma}$ is positive semi-definite, thus all eigenvalues are non-negative.

PCA

Let's start with $K = 1$,

$$\mathbf{u}_1 = \arg \max_{\|\mathbf{u}\|_2^2=1} \mathbf{u}^T \hat{\Sigma} \mathbf{u}$$

Lagrangian:

$$L(\mathbf{u}, \lambda) = \mathbf{u}^T \hat{\Sigma} \mathbf{u} + \lambda(1 - \|\mathbf{u}\|_2^2)$$

Setting the gradient to zero: $2\hat{\Sigma}\mathbf{u} - 2\lambda\mathbf{u} = 0 \Rightarrow \hat{\Sigma}\mathbf{u} = \lambda\mathbf{u}$,

i.e., \mathbf{u} is an eigenvector of $\hat{\Sigma}$. Plugging above:

$$\mathbf{u}^T \hat{\Sigma} \mathbf{u} = \mathbf{u}^T \mathbf{u} \lambda = \|\mathbf{u}\|_2^2 \lambda = \lambda$$

Conclusion: \mathbf{u}_1 is the eigenvector of the largest eigenvalue of $\hat{\Sigma}$.

Notice: $\hat{\Sigma}$ is positive semi-definite, thus all eigenvalues are non-negative.

Easy to extend to $K > 1$.

Interlude: PCA

$$\arg \max_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \sum_{j=1}^K \mathbf{u}_j^T \hat{\Sigma} \mathbf{u}_j = \text{eigenvectors of the top } K \text{ eigenvalues of } \hat{\Sigma}$$

$$\text{Recall that } \hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T = \frac{1}{N} \mathbf{X}^T \mathbf{X}$$

Interlude: PCA

$$\arg \max_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \sum_{j=1}^K \mathbf{u}_j^T \hat{\Sigma} \mathbf{u}_j = \text{eigenvectors of the top } K \text{ eigenvalues of } \hat{\Sigma}$$

Recall that $\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T = \frac{1}{N} \mathbf{X}^T \mathbf{X}$

Notice that the eigenvalues of $\hat{\Sigma}$ are the **singular values** (SV) of \mathbf{X}/\sqrt{N} .

PCA corresponds to SVD (SV decomposition) of the scaled data matrix \mathbf{X}

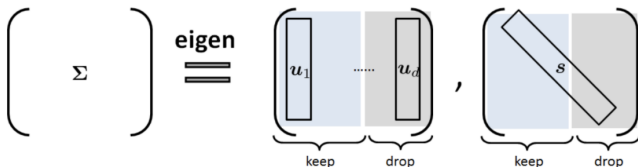
Interlude: PCA

$$\arg \max_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \sum_{j=1}^K \mathbf{u}_j^T \hat{\Sigma} \mathbf{u}_j = \text{eigenvectors of the top } K \text{ eigenvalues of } \hat{\Sigma}$$

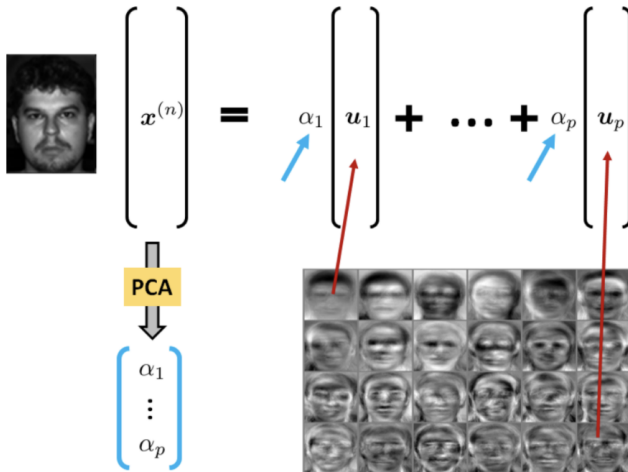
Recall that $\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T = \frac{1}{N} \mathbf{X}^T \mathbf{X}$

Notice that the eigenvalues of $\hat{\Sigma}$ are the **singular values** (SV) of \mathbf{X}/\sqrt{N} .

PCA corresponds to SVD (SV decomposition) of the scaled data matrix \mathbf{X}



PCA: EigenFaces



Back to the Linear Auto-Encoder

Let $\mathbf{X} \in \mathbb{R}^{N \times D}$ be the data matrix ($N > D$),

Back to the Linear Auto-Encoder

Let $\mathbf{X} \in \mathbb{R}^{N \times D}$ be the data matrix ($N > D$),

Assume $\mathbf{W} \in \mathbb{R}^{K \times D}$ with $K < D$ (no biases, assuming \mathbf{X} is centred)

Back to the Linear Auto-Encoder

Let $\mathbf{X} \in \mathbb{R}^{N \times D}$ be the data matrix ($N > D$),

Assume $\mathbf{W} \in \mathbb{R}^{K \times D}$ with $K < D$ (no biases, assuming \mathbf{X} is centred)

We want to minimize, with g identity, thus $\hat{\mathbf{x}}^{(i)} = \mathbf{W}^\top \mathbf{W} \mathbf{x}^{(i)}$,

$$\sum_{i=1}^N \|\mathbf{x}^{(i)} - \hat{\mathbf{x}}^{(i)}\|_2^2 = \|\mathbf{X} - \hat{\mathbf{X}}\|_F^2 = \|\mathbf{X} - \mathbf{X} \mathbf{W}^\top \mathbf{W}\|_F^2$$

where $\|\cdot\|_F^2$ is the Frobenius matrix norm and $\mathbf{W}^\top \mathbf{W}$ has rank K .

Back to the Linear Auto-Encoder

Let $\mathbf{X} \in \mathbb{R}^{N \times D}$ be the data matrix ($N > D$),

Assume $\mathbf{W} \in \mathbb{R}^{K \times D}$ with $K < D$ (no biases, assuming \mathbf{X} is centred)

We want to minimize, with g identity, thus $\hat{\mathbf{x}}^{(i)} = \mathbf{W}^\top \mathbf{W} \mathbf{x}^{(i)}$,

$$\sum_{i=1}^N \|\mathbf{x}^{(i)} - \hat{\mathbf{x}}^{(i)}\|_2^2 = \|\mathbf{X} - \hat{\mathbf{X}}\|_F^2 = \|\mathbf{X} - \mathbf{X} \mathbf{W}^\top \mathbf{W}\|_F^2$$

where $\|\cdot\|_F^2$ is the Frobenius matrix norm and $\mathbf{W}^\top \mathbf{W}$ has rank K .

From the Eckart-Young theorem, the minimizer is **truncated SVD** of \mathbf{X}^\top :

$$\hat{\mathbf{X}}^\top = \mathbf{U}_K \Sigma_K \mathbf{V}_K^\top,$$

where Σ_K is a diagonal matrix containing the top K singular values of \mathbf{X}^\top , and the columns of \mathbf{U}_K are the corresponding left singular vectors.

Back to the Linear Auto-Encoder

Let $\mathbf{X} \in \mathbb{R}^{N \times D}$ be the data matrix ($N > D$),

Assume $\mathbf{W} \in \mathbb{R}^{K \times D}$ with $K < D$ (no biases, assuming \mathbf{X} is centred)

We want to minimize, with g identity, thus $\hat{\mathbf{x}}^{(i)} = \mathbf{W}^\top \mathbf{W} \mathbf{x}^{(i)}$,

$$\sum_{i=1}^N \|\mathbf{x}^{(i)} - \hat{\mathbf{x}}^{(i)}\|_2^2 = \|\mathbf{X} - \hat{\mathbf{X}}\|_F^2 = \|\mathbf{X} - \mathbf{XW}^\top \mathbf{W}\|_F^2$$

where $\|\cdot\|_F^2$ is the Frobenius matrix norm and $\mathbf{W}^\top \mathbf{W}$ has rank K .

From the Eckart-Young theorem, the minimizer is **truncated SVD** of \mathbf{X}^\top :

$$\hat{\mathbf{X}}^\top = \mathbf{U}_K \Sigma_K \mathbf{V}_K^\top,$$

where Σ_K is a diagonal matrix containing the top K singular values of \mathbf{X}^\top , and the columns of \mathbf{U}_K are the corresponding left singular vectors.

The solution is $\mathbf{W} = \mathbf{U}_K^\top$, which gives as desired:

$$\hat{\mathbf{X}}^\top = \mathbf{W}^\top \mathbf{W} \mathbf{X}^\top = \mathbf{U}_K \mathbf{U}_K^\top \mathbf{U} \Sigma \mathbf{V}^\top = \mathbf{U}_K \Sigma_K \mathbf{V}_K^\top.$$

Back to the Linear Auto-Encoder

Let $\mathbf{X} \in \mathbb{R}^{N \times D}$ be the data matrix ($N > D$),

Assume $\mathbf{W} \in \mathbb{R}^{K \times D}$ with $K < D$ (no biases, assuming \mathbf{X} is centred)

We want to minimize, with g identity, thus $\hat{\mathbf{x}}^{(i)} = \mathbf{W}^\top \mathbf{W} \mathbf{x}^{(i)}$,

$$\sum_{i=1}^N \|\mathbf{x}^{(i)} - \hat{\mathbf{x}}^{(i)}\|_2^2 = \|\mathbf{X} - \hat{\mathbf{X}}\|_F^2 = \|\mathbf{X} - \mathbf{X} \mathbf{W}^\top \mathbf{W}\|_F^2$$

where $\|\cdot\|_F^2$ is the Frobenius matrix norm and $\mathbf{W}^\top \mathbf{W}$ has rank K .

From the Eckart-Young theorem, the minimizer is **truncated SVD** of \mathbf{X}^\top :

$$\hat{\mathbf{X}}^\top = \mathbf{U}_K \Sigma_K \mathbf{V}_K^\top,$$

where Σ_K is a diagonal matrix containing the top K singular values of \mathbf{X}^\top , and the columns of \mathbf{U}_K are the corresponding left singular vectors.

The solution is $\mathbf{W} = \mathbf{U}_K^\top$, which gives as desired:

$$\hat{\mathbf{X}}^\top = \mathbf{W}^\top \mathbf{W} \mathbf{X}^\top = \mathbf{U}_K \mathbf{U}_K^\top \mathbf{U} \Sigma \mathbf{V}^\top = \mathbf{U}_K \Sigma_K \mathbf{V}_K^\top.$$

Conclusion: the optimal linear auto-encoder coincides with PCA.

Auto-Encoders

PCA fits a **linear manifold** (affine space) to the data

By using **non-linear** activations, we obtain more sophisticated codes (i.e. representations).

We need some sort of regularization to:

- encourage a smooth representation (small perturbations of the input will lead to similar codes)
- avoid overfitting to the training data

Some Variants of Auto-Encoders

- **Sparse auto-encoders:** use many hidden units, but add a ℓ_1 regularization term to encourage **sparse representations** of the input
- **Denoising auto-encoders:** regularize by adding noise to the input; the goal is to learn a **smooth representation function** that allows to output the denoised input (inspired by image denoising)
- **Stacked auto-encoders:** several auto-encoders on top of each other
- **Variational auto-encoders:** a generative probabilistic model that minimizes a variational bound (this will be covered in another lecture!)

Regularized Auto-Encoders

- To regularize auto-encoders, **regularization** may be added to the loss
- The goal is then to minimize $L(\hat{x}; x) + \Omega(h, x)$
- For example:
 - regularizing the code $\Omega(h, x) = \lambda \|h\|^2$
 - regularizing the derivatives $\Omega(h, x) = \lambda \sum_i \|\nabla_x h_i\|^2$
- The encoder and decoder parameters may be shared or not.

Sparse Auto-Encoders

- Most auto-encoders learn low-dimensional codes, e.g., they reduce input dimensionality (bottleneck shape $K < D$).
- But one exception are **sparse auto-encoders**:
 - Sparse auto-encoders incorporate a sparsity penalty $\Omega(\mathbf{h})$ on the code layer, e.g., $\Omega(\mathbf{h}) = \lambda \|\mathbf{h}\|_1$
 - Typically the number of hidden units is large, e.g., larger than the input dimension
 - The sparsity penalty encourages sparse codes, where most hidden units are inactive.

Stochastic Auto-Encoders

- In this case, the encoder and decoder are not deterministic functions, but involve some noise injection
- We have a distribution $p_{\text{encoder}}(\mathbf{h} \mid \mathbf{x})$ for the encoder and a distribution $p_{\text{decoder}}(\mathbf{x} \mid \mathbf{h})$ for the decoder
- The auto-encoder can be trained to minimize

$$-\log p_{\text{decoder}}(\mathbf{x} \mid \mathbf{h}).$$

Denoising Auto-Encoders

- Use a perturbed version of the input, $\tilde{\mathbf{x}} = \mathbf{x} + \mathbf{n}$, where \mathbf{n} is random noise (e.g. Gaussian noise $\mathbf{n} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$)
- Instead of minimizing $\frac{1}{2} \|\hat{\mathbf{x}} - \mathbf{x}\|^2$, minimize $\frac{1}{2} \|\hat{\mathbf{x}} - \tilde{\mathbf{x}}\|^2$
- This is a form of implicit regularization that ensures **smoothness**: it forces the system to represent well not only the data points, but also their perturbations

Denoising Auto-Encoders

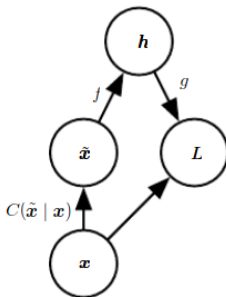


Figure 14.3: The computational graph of the cost function for a denoising autoencoder, which is trained to reconstruct the clean data point \mathbf{x} from its corrupted version $\tilde{\mathbf{x}}$. This is accomplished by minimizing the loss $L = -\log p_{\text{decoder}}(\mathbf{x} | \mathbf{h} = f(\tilde{\mathbf{x}}))$, where $\tilde{\mathbf{x}}$ is a corrupted version of the data example \mathbf{x} , obtained through a given corruption process $C(\tilde{\mathbf{x}} | \mathbf{x})$. Typically the distribution p_{decoder} is a factorial distribution whose mean parameters are emitted by a feedforward network g .

(From Goodfellow et al.'s book.)

Denoising Auto-Encoders

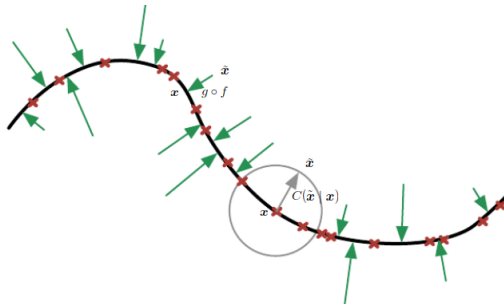


Figure 14.4: A denoising autoencoder is trained to map a corrupted data point $\tilde{\mathbf{x}}$ back to the original data point \mathbf{x} . We illustrate training examples \mathbf{x} as red crosses lying near a low-dimensional manifold, illustrated with the bold black line. We illustrate the corruption process $C(\tilde{\mathbf{x}} | \mathbf{x})$ with a gray circle of equiprobable corruptions. A gray arrow demonstrates how one training example is transformed into one sample from this corruption process. When the denoising autoencoder is trained to minimize the average of squared errors $\|g(f(\tilde{\mathbf{x}})) - \mathbf{x}\|^2$, the reconstruction $g(f(\tilde{\mathbf{x}}))$ estimates $\mathbb{E}_{\mathbf{x}, \tilde{\mathbf{x}} \sim p_{\text{data}}(\mathbf{x}) C(\tilde{\mathbf{x}} | \mathbf{x})} [\mathbf{x} | \tilde{\mathbf{x}}]$. The vector $g(f(\tilde{\mathbf{x}})) - \tilde{\mathbf{x}}$ points approximately toward the nearest point on the manifold, since $g(f(\tilde{\mathbf{x}}))$ estimates the center of mass of the clean points \mathbf{x} that could have given rise to $\tilde{\mathbf{x}}$. The autoencoder thus learns a vector field $g(f(\mathbf{x})) - \mathbf{x}$ indicated by the green arrows. This vector field estimates the score $\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$ up to a multiplicative factor that is the average root mean square reconstruction error.

Why Do We Use Auto-Encoders?

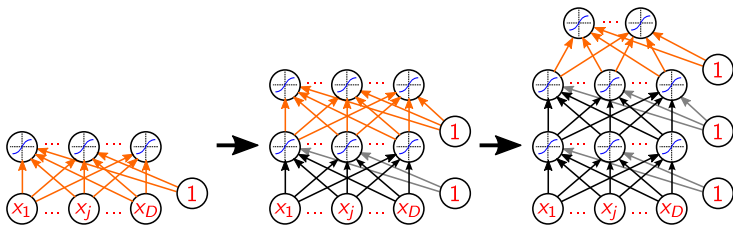
Historically, training deep neural networks was hard

One of the initial successful uses of auto-encoders was for **unsupervised pre-training** (Erhan et al., 2010).

Unsupervised Pre-Training

A greedy, layer-wise procedure:

- train one layer at a time, from first to last, with unsupervised criterion (e.g. an auto-encoder)
- fix the parameters of previous hidden layers
- previous layers viewed as feature extraction



Pre-training initializes the parameters in a region such that the near local optima overfit less the data.

Fine-Tuning

Once all layers are pre-trained:

- add output layer
- train the whole network using supervised learning

Supervised learning is performed as in a regular feed-forward network:

- forward propagation, backpropagation and update
- all parameters are “tuned” for the supervised task at hand
- representation is adjusted to be more discriminative

Other Applications of Auto-Encoders

- Dimensionality reduction
- Information retrieval and semantic hashing (via binarizing the codes)
- Conversion of discrete inputs to low-dimensional continuous space

Outline

① Representation Learning

Hierarchical Compositionality

Distributed Representations

Auto-Encoders

Word Embeddings

② Conclusions

Word Representations

We'll focus now on recent methods for learning representations of **words in natural language**

Also called **word embeddings**

This has been an extremely successful application of representation learning

It's still a very active area of research!

Distributional Similarity

Key idea: represent a word by means of its neighbors

- “*You shall know a word by the company it keeps*” (J. R. Firth, 1957)
- One of the most successful ideas of modern statistical NLP!

For example:

- Adjectives are normally surrounded by nouns
- Words like *book*, *newspaper*, *article*, are commonly surrounded by *reading*, *read*, *writes*, but not by *flying*, *eating*, *sleeping*

Word Embeddings

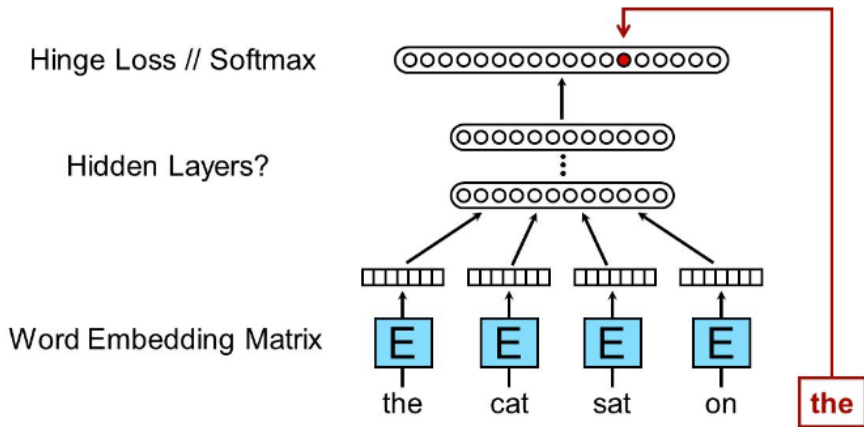
How do we obtain lower dimensional vector representations of words?

Two possible methods:

- Factorization of a co-occurrence word/context matrix (latent semantic analysis, etc.)
- Directly learn low-dimensional vectors by training a network to *predict* the context of a given word

We'll focus on the latter, incarnated in the **word2vec** toolkit (Mikolov et al., 2013), which follows previous ideas of Bengio et al. (2003) and Collobert et al. (2011).

Neural Language Model (Bengio et al., 2003)



(Image credits: Quoc Le)

Neural Language Model (Bengio et al., 2003)

- Each word is associated with a continuous vector (a **word embedding**)
- Given the **context** (previous K words), predict the next word
- This is done by concatenating the word embeddings in the context window, then propagating them through a feedforward neural network
- The output layer is a gigantic softmax that assigns a probability value to each word in the vocabulary

Variants of this model achieved better accuracy than smoothed K -th order Markov models

As a by-product: **word embeddings!**

The embedding matrix is a lookup table that assigns a continuous vector to every word in the vocabulary.

Neural Language Model

In this class, we are not concerned with language modeling (the actual task), but rather about the **quality of the embeddings** (the representations we learn for that task).

Some Insights

If we don't care about language modeling as a task:

- ① We don't need to have a “left-to-right model” where we try to predict the next word given the context
- ② We don't need to predict the probability of every word, we might just make sure that the true word is more likely than a random word

These insights underlie the word2vec model of Mikolov et al. (2013).

Word2Vec (Mikolov et al., 2013)

Considers a context window **around** each word in the sentence.

Word2vec comes with two variants:

- **Skip-gram**: predict surrounding context words in a window of length m of every word
- **Continuous bag-of-words** (CBOW): predict the central word from the context

We'll focus on the skip-gram model (more widely used).

Skip-Gram

Goal: maximize the log probability of any context word given the current center word:

$$J(\Theta) = \frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log p_{\Theta}(x_{t+j} \mid x_t)$$

There are two sets of parameters $\Theta = (\mathbf{u}, \mathbf{v})$:

- Embeddings \mathbf{u}_o for each word o appearing as the center word
- Embeddings \mathbf{v}_c for each word c appearing in the context of another word

Define a **log-bilinear model**: $p_{\Theta}(x_{t+j} = c \mid x_t = o) \propto \exp(\mathbf{u}_o \cdot \mathbf{v}_c)$

Every word gets two vectors!

In the end, we use the \mathbf{u} vectors as the word embeddings and discard the \mathbf{v} vectors

The Large Vocabulary Problem

Recall that we have

$$p_{\Theta}(x_{t+j} = c \mid x_t = o) = \frac{\exp(\mathbf{u}_o \cdot \mathbf{v}_c)}{\sum'_c \exp(\mathbf{u}_o \cdot \mathbf{v}'_c)}$$

This objective requires a softmax over the entire vocabulary

Unfortunately, with large vocabularies this leads to very slow training :(

Workarounds:

- Stochastic sampling
- Noise contrastive estimation
- Negative sampling

More details in these notes: <https://arxiv.org/pdf/1410.8251.pdf>

We'll focus on **negative sampling**.

Negative Sampling

Key idea:

- replace the gigantic softmax by **binary logistic regressions** for a true pair (center word and word in its context window) and a couple of random pairs (the center word with a random word):

$$J_t(\Theta) = \log \sigma(\mathbf{u}_o \cdot \mathbf{v}_c) + \sum_{i=1}^k \log \sigma(-\mathbf{u}_o \cdot \mathbf{v}_{j_i}), \quad j_i \sim P(x)$$

- Several strategies for the sampling distribution $P(x)$ (uniform, unigram frequency, etc.)

Negative sampling is a simple form of unsupervised pre-training.

Linear Relationships

- These representations are very good at encoding dimensions of similarity!
- **Word analogies** can be solved quite well just by doing vector subtraction in the embedding space
- Syntactically:

$$\mathbf{x}_{\text{apple}} - \mathbf{x}_{\text{apples}} \approx \mathbf{x}_{\text{car}} - \mathbf{x}_{\text{cars}} \approx \mathbf{x}_{\text{family}} - \mathbf{x}_{\text{families}}$$

- Semantically:

$$\begin{aligned}\mathbf{x}_{\text{shirt}} - \mathbf{x}_{\text{clothing}} &\approx \mathbf{x}_{\text{chair}} - \mathbf{x}_{\text{furniture}} \\ \mathbf{x}_{\text{king}} - \mathbf{x}_{\text{man}} &\approx \mathbf{x}_{\text{queen}} - \mathbf{x}_{\text{woman}}\end{aligned}$$

Visualization

Typical word embedding dimensions are on the hundreds (e.g. 300)

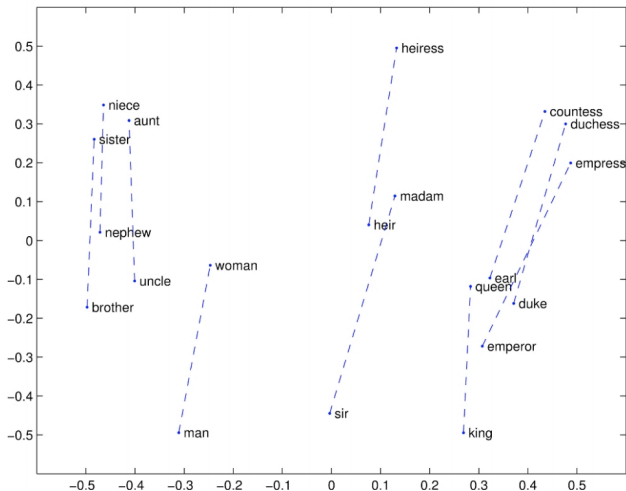
How can we visualize these embeddings?

Simple way: project them in 2D with something like PCA!

Most used: t -distributed stochastic neighbor embedding (t-SNE, Maaten and Hinton 2008)

<https://lvdmaaten.github.io/tsne>

Word Analogies (Mikolov et al., 2013)



(Slide credit to Richard Socher)

Other Methods for Obtaining Word Embeddings

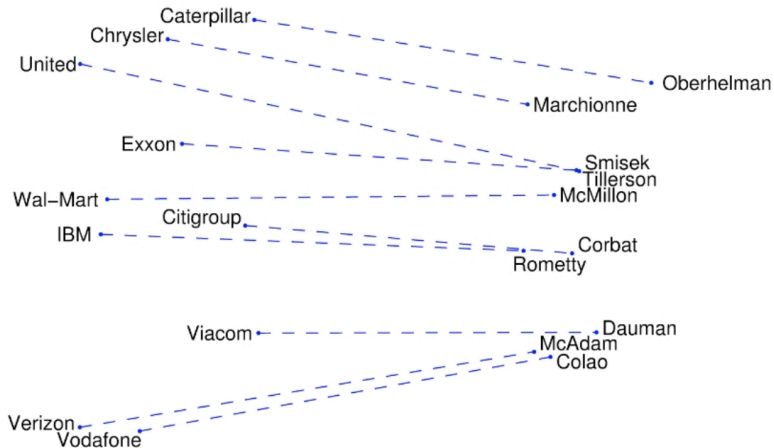
GloVe: Global Vectors for Word Representation (Pennington et al., 2014)

- <https://nlp.stanford.edu/projects/glove>
- Training is performed on aggregated global word-word co-occurrence statistics from a corpus

fastText (Bojanowski et al., 2016): embeds also character n -grams for generating embeddings for out-of-vocabulary words

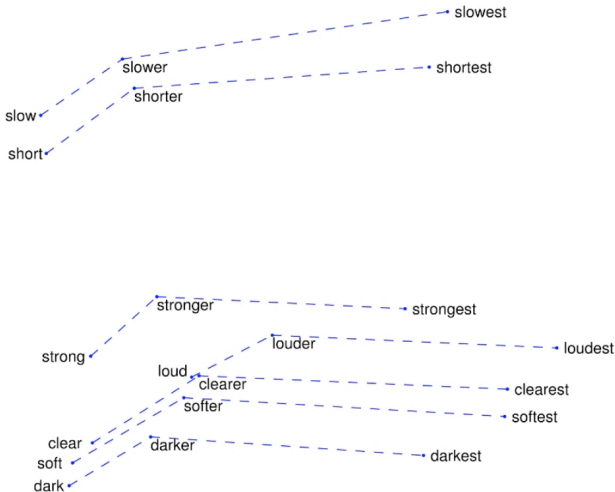
- <https://fasttext.cc> (from FAIR)
- open-source, free, lightweight library that allows users to learn text representations and text classifiers
- contains multi-lingual word vectors for 157 different languages

GloVe Visualizations: Company \rightarrow CEO



(Slide credit to Richard Socher)

GloVe Visualizations: Superlatives



(Slide credit to Richard Socher)

Word Embeddings: Some Open Problems

- Can we have word embeddings for multiple languages in the same space?
- How to capture polysemy?
- These word embeddings are static, can we compute embeddings on-the-fly depending on the context?

Cross-Lingual Word Embeddings

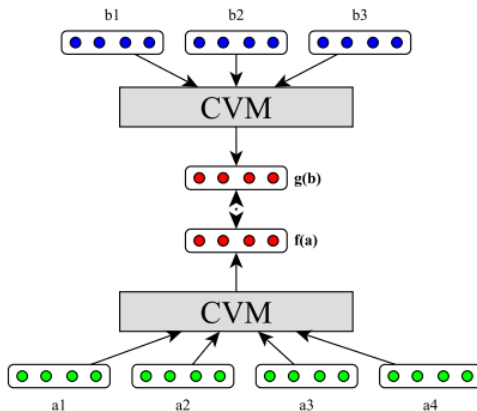


Figure 1: Model with parallel input sentences a and b . The model minimises the distance between the sentence level encoding of the bitext. Any composition functions (CVM) can be used to generate the compositional sentence level representations.

(From Hermann and Blunsom (2014).)

Cross-Lingual Word Embeddings

Key idea:

- use a corpus of parallel sentences in two languages
- define a composition function to obtain a sentence representation given word embeddings
- apply a loss function that encourages the sentence representations in the two languages to be similar
- negative sampling works here too: true pair vs fake pair.

Cross-Lingual Word Embeddings

Other approaches:

- Define a bilingual dictionary and apply canonical correlation analysis (Faruqui and Dyer, 2014)
- Task-specific embeddings with convex optimization (Ferreira et al., 2016)
- Learn the two embeddings separately, and then apply a linear transformation to put them in a shared space (Artetxe et al., 2017)
- Adversarial training (Lample et al., 2018)

This is a very active area of research!

Contextual Embeddings

Words can have different meanings, depending on which context they appear in.

In 2018, a model called ELMo learned **context-dependent embeddings** and achieved impressive results on 6 NLP downstream tasks (Peters et al., 2018)

Key idea:

- Pre-train a BiLSTM language model on a large dataset (we'll see in a later class what this is)
- Save **all** the encoder parameters at all layers, not only the embeddings
- Then, for your downstream task, tune a scalar parameter for each layer, and pass **the entire sentence** through this encoder.

BERT, GPT, etc.

Some time later, a Transformer-based model (BERT) achieved even better performance:

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova

Google AI Language

`{jacobdevlin, mingweichang, kentonl, kristout}@google.com`

Huge improvements in multiple NLP tasks!

(Trained on 64 TPU chips!!)

Other related models include GPT-2, GPT-3, etc.

This will be covered in a later lecture!

Outline

① Representation Learning

Hierarchical Compositionality

Distributed Representations

Auto-Encoders

Word Embeddings

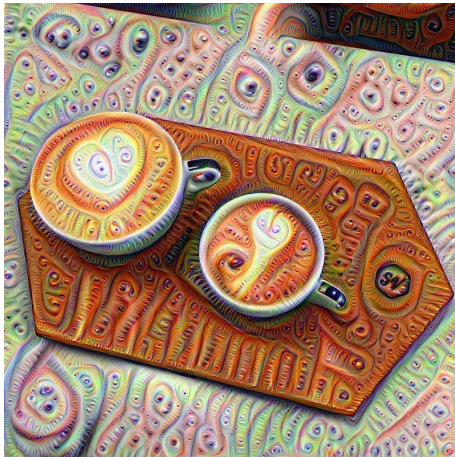
② Conclusions

Conclusions

- Neural nets learn internal representations that can be transferred across tasks
- Distributed representations are exponentially more compact and allow generalizing to unseen objects
- Deeper neural nets exhibit hierarchical compositionality: upper level layers learn more abstract/semantic representations than bottom level layers
- Auto-encoders are an effective means for learning representations
- Word embeddings are continuous representations of words that are extremely useful in NLP

Thank you!

Questions?



References I

- Artetxe, M., Labaka, G., and Agirre, E. (2017). Learning bilingual word embeddings with (almost) no bilingual data. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 451–462.
- Bengio, Y., Ducharme, R., Vincent, P., and Jauvin, C. (2003). A neural probabilistic language model. *Journal of Machine Learning Research*, 3(Feb):1137–1155.
- Bojanowski, P., Grave, E., Joulin, A., and Mikolov, T. (2016). Enriching word vectors with subword information. *arXiv preprint arXiv:1607.04606*.
- Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., and Kuksa, P. (2011). Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12(Aug):2493–2537.
- Erhan, D., Bengio, Y., Courville, A., Manzagol, P.-A., Vincent, P., and Bengio, S. (2010). Why does unsupervised pre-training help deep learning? *Journal of Machine Learning Research*, 11(Feb):625–660.
- Faruqui, M. and Dyer, C. (2014). Improving vector space word representations using multilingual correlation. In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 462–471.
- Ferreira, D., Almeida, M. S. C., and Martins, A. F. T. (2016). Jointly Learning to Embed and Predict with Multiple Languages. In *Proc. of the Annual Meeting of the Association for Computational Linguistics*.
- Hermann, K. M. and Blunsom, P. (2014). Multilingual Models for Compositional Distributional Semantics. In *Proc. of the Annual Meeting of the Association for Computational Linguistics*.
- Hinton, G. E. (1984). Distributed representations.
- Lample, G., Ott, M., Conneau, A., Denoyer, L., and Ranzato, M. (2018). Phrase-based & neural unsupervised machine translation. *arXiv preprint arXiv:1804.07755*.
- Maaten, L. v. d. and Hinton, G. (2008). Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119.
- Pennington, J., Socher, R., and Manning, C. D. (2014). Glove: Global Vectors for Word Representation. *Proceedings of the Empirical Methods in Natural Language Processing (EMNLP 2014)*, 12:1532–1543.
- Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., and Zettlemoyer, L. (2018). Deep contextualized word representations. *arXiv preprint arXiv:1802.05365*.
- Zeiler, M. D. and Fergus, R. (2013). Stochastic pooling for regularization of deep convolutional neural networks. *arXiv preprint arXiv:1301.3557*.