MATH63800 Mini-Project 1 Feature Extraction and Transfer Learning on Fashion-MNIST

> **Jason WU**, Peng XU, Nayeon LEE 08.Mar.2018

Introduction: Fashion-MNIST Dataset

- 60,000 training examples and a 10,000 testing examples
- Each example is a 28x28 grayscale image
- 10 classes

index	0	1	2	3	4	5	6	7	8	9
Туре	T-shirt/top	Trouser	Pullover	Dress	Coat	Sandal	Shirt	Sneaker	Bag	Ankle boot

• Zalando et al. intend Fashion-MNIST to serve as a direct drop-in replacement for the original MNIST dataset for benchmarking machine learning algorithms.

Why Fashion-MNIST?

Quoted from their website:

- MNIST is too easy. Convolutional nets can achieve 99.7% on MNIST. Classic machine learning algorithms can also achieve 97% easily. Most pairs of MNIST digits can be distinguished pretty well by just one pixel.
- MNIST is overused. In this April 2017 Twitter thread, Google Brain research scientist and deep learning expert Ian Goodfellow calls for people to move away from MNIST.
- MNIST can not represent modern CV tasks, as noted in this April 2017 Twitter thread, deep learning expert/Keras author François Chollet.

Introduction: Fashion-MNIST Dataset



▋▝▋▓▌▌▋▋▆▓▋▌▓▋▆▓▋▆▋▋▓▊▋▆▋▋▆▋▋▓▋▋▓▆▓▋ ▆▆ড়ড়ড়ॺ⋨ऒऒ▆▆▓▓▆▅▋▆▆▓▋▆▆▋▋▆▆₽₩
₹ ₽ ₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽
Ţ <u>, , , , , , , , , , , , , , , , , , , </u>
~ De a G ~ A D D D ~ Z A J ~ A
≁ℳ <i>⋥⋘⋬⋐⋈∼≁∽∽⋏⋐⋥</i> ⋴⋩⋰ਔ⋐ ℒ⋏ ₰₶∕⋏⋷ख़∝≖⋇∂ Т₿₿₿₿₱₿⋒₽₿Ⅱ₲⋒⋒常 Т₿ ⋒₽₤⋒₨⋒⋒⋒₿₿⋒₿₿
▝▋▓▓▌▓▓▟▆▆▆▆▆▆▆▆▆▆▆▆▆▆▆▋▓▋▋ ▆▓▓▌▓▓▟▆▆▆▆▆▆▆▆▆▆▆▆▆▆▆▆▆▆₽₽₽₽₽₽₽₽₽₽₽₽₽₽₽
a a da a a a a a a a a a a a a a a a a
A A A A A A A A A A A A A A A A A A A
Ŀĸĸਙ&Ŀ& Ŗ ≈₽₽₽₽₽ Ę <u>Ę</u> <u>₹₽₽₽</u> <u>₽₽₽₽₽₽₽₽</u> <u>® ĸ₽₽₽₽₽</u> ₽₽₽₽₽₽₽₽₽₽₽₽₽ <u>₽</u> ₽₽₽ <u>₽</u> ₽₽₽₽₽₽₽

Material: https://github.com/zalandoresearch/fashion-mnist

How to import?

- Loading data with Python (requires NumPy)
 - Use utils/mnist_reader from <u>https://github.com/zalandoresearch/fashion-mnist</u>

```
import mnist_reader
X_train, y_train = mnist_reader.load_mnist('data/fashion', kind='train')
X_test, y_test = mnist_reader.load_mnist('data/fashion', kind='t10k')
```

- Loading data with Tensorflow
 - Make sure you have downloaded the data and placed it in data/fashion.
 Otherwise, Tensorflow will download and use the original MNIST.

```
from tensorflow.examples.tutorials.mnist import input_data
data = input_data.read_data_sets('data/fashion')
```

```
data.train.next_batch(BATCH_SIZE)
```

Feature Extraction

- We compared three different feature representation:
 - Raw pixel features
 - ScatNet features
 - Pretrained ResNet18 last-layer features















Feature Extraction(1): ScatNet

- The maximum scale of the transform:]=3
- The maximum scattering order: M=2
- The number of different orientations: L=1

The dimension of the final features is 176



https://arxiv.org/pdf/1203.1513.pdf

Feature Extraction(2): ResNet

- Used pretrained 18 layers Residual Network from ImageNet
- We take the hidden representation right before the last fully-connected layer, which has the dimension of 512



https://arxiv.org/abs/1512.03385

Data Visualization

- Then, we visualized three different feature representation by the following 4 different dimension reduction methods:
 - Principal Component Analysis (PCA)
 - Locally Linear Embedding (LLE)
 - t-Distributed Stochastic Neighbor Embedding (t-SNE)
 - Uniform Manifold Approximation and Projection (UMAP)

index	0	1	2	3	4	5	6	7	8	9
Type	T-shirt/top	Trouser	Pullover	Dress	Coat	Sandal	Shirt	Sneaker	Bag	Ankle boot
Color	'magenta'	'cyan'	'red'	'green'	'blue'	'black'	'yellow'	'orange'	'navy'	'gray'

Data Visualization

- Then, we visualized three different feature representation by the following 4 different dimension reduction methods:
 - Principal Component Analysis (PCA)
 - Locally Linear Embedding (LLE)
 - t-Distributed Stochastic Neighbor Embedding (t-SNE)
 - Uniform Manifold Approximation and Projection (UMAP)



Data Visualization (1): PCA



- Normalization, Covariance Matrix, SVD, Project to top K eigen-vectors
- Linear dimension reduction methods:
 - not that obviously difference between labels

Data Visualization (2): LLE



- 1. Compute the neighbors of each data point, \vec{X}_{i} .
- 2. Compute the weights W_{ij} that best reconstruct each data point \vec{X}_i from its neighbors, minimizing the cost in Equation (1) by constrained linear fits.
- 3. Compute the vectors \vec{Y}_i best reconstructed by the weights W_{ij} , minimizing the quadratic form in Equation (2) by its bottom nonzero eigenvectors.



http://www.robots.ox.ac.uk/~az/lectures/ml/lle.pdf

Data Visualization (2): LLE



Figure 1. A nonlinear S-manifold consisting of 1000 data points



a) the result obtained by LLE



b) the result obtained by PCA



Figure 3. Embeddings of the 2-dimensional S-manifold, computed for different choices of the number of the nearest neighbours *k* by LLE

https://pdfs.semanticscholar.org/6adc/19cf4404b9f1224a1a027022e40ac77218f5.pdf

Data Visualization (2): LLE



Non-linear dimension reduction that is good at capture "streamline" structure

Data Visualization (3): t-SNE

- Use Gaussian pdf to approximate the high dimension distribution
- Use t distribution for low dimension distribution
- Use KL Divergence as cost function for gradient descent

$$p_{ij} = \frac{\exp\frac{\left(-\|x_i - x_j\|^2\right)}{2\sigma^2}}{\sum\frac{\exp(-\|x_k - x_l\|^2)}{2\sigma^2}} \qquad \qquad q_{ij} = \frac{\left(1 + \|y_i - y_j\|^2\right)^{-1}}{\sum\left(1 + \|y_k - y_l\|^2\right)^{-1}}$$

$$C = KL(P||Q) = \sum_{i} \sum_{j} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$
$$\frac{\delta C}{\delta y_i} = 4\sum_{j} (p_{ij} - q_{ij}) (y_i - y_j)(1 + ||y_i - y_j||^2)^{-1}$$

http://www.jmlr.org/papers/volume9/vandermaaten08a/vandermaaten08a.pdf

Data Visualization (3): t-SNE



• Block-like visualization due to the gaussian approximation

Data Visualization (4): UMAP

- The algorithm is founded on three assumptions about the data
 - The Riemannian metric is locally constant (or can be approximated);
 - The data is uniformly distributed on Riemannian manifold;
 - The manifold is locally connected.

https://arxiv.org/pdf/1802.03426.pdf



Data Visualization (4): UMAP

https://github.com/lmcinnes/umap



 Much more Faster in training process, which implies it can handle large datasets and high dimensional data

Any News from Visualization?

- Is there different patterns between different visualization methods?
- Is there clear separation of different classes?
- Is there any groups that tend to cluster together?
- Let's look closer!

PCA



LLE



t-SNE



UMAP



Sneaker, Sandal, Ankle boot



PCA



LLE



t-SNE



UMAP



Trouser



PCA



LLE



t-SNE



UMAP





Bag

PCA



LLE



t-SNE



UMAP



T-Shirt, Pullover, Dress, Coat, Shirt



- Logistic Regression
- Linear Discriminant Analysis
- Support Vector Machine
- Random Forest

...

Logistic Regression



- Linear Discriminant Analysis
 - maximize between class covariance
 - minimize within class covariance

-2

2

6

$$S = \frac{\sigma_{\text{between}}^2}{\sigma_{\text{within}}^2} = \frac{(\vec{w} \cdot \vec{\mu}_1 - \vec{w} \cdot \vec{\mu}_0)^2}{\vec{w}^T \Sigma_1 \vec{w} + \vec{w}^T \Sigma_0 \vec{w}} = \frac{(\vec{w} \cdot (\vec{\mu}_1 - \vec{\mu}_0))^2}{\vec{w}^T (\Sigma_0 + \Sigma_1) \vec{w}}$$

-2

Linear Support Vector Machine

 Hard-margin

$$ec{w}\cdotec{x}_i-b\geq 1,$$
 if $y_i=1$

or

$$ec{w}\cdotec{x}_i-b\leq -1,$$
 if $y_i=-1.$

• Soft-margin

$$\left[rac{1}{n}\sum_{i=1}^n \max\left(0,1-y_i(ec{w}\cdotec{x}_i-b)
ight)
ight]+\lambda \|ec{w}\|^2,$$



- Random Forest
- An ensemble learning method that construct multiple decision trees
- Bagging (Bootstrap aggregating)

For *b* = 1, ..., *B*:

- 1. Sample, with replacement, *n* training examples from *X*, *Y*; call these X_b , Y_b .
- 2. Train a classification or regression tree f_b on X_b , Y_b .



Simple Classification Results

Raw Features 784	ScatNet Features 176	ResNet Features 512
84.60%	80.83%	79.03%
81.51%	79.96%	77.08%
78.84%	82.60%	73.63%
87.12%	83.06%	59.20%
83.01%	81.61%	72 24%
	Raw Features 784 84.60% 81.51% 78.84% 87.12% 83.01%	Raw Features 784 ScatNet Features 176 84.60% 80.83% 81.51% 79.96% 78.84% 82.60% 87.12% 83.06% 83.01% 81.61%

Simple Classification Results

Models	Raw Features 784	ScatNet Features 176	ResNet Features 512	Average
Logistic Regression	84.60%	80.83%	79.03%	81.49%
Linear Discriminant Analysis (LDA)	81.51%	79.96%	77.08%	79.52%
Support Vector Machine (SVM)	78.84%	82.60%	73.63%	78.36%
Random Forest	87.12%	83.06%	59.20%	76.46%

Simple Classification Results

Models	Raw Features 784	ScatNet Features 176	ResNet Features 512
Logistic Regression	84.60%	80.83%	79.03%
Linear Discriminant Analysis (LDA)	81.51%	79.96%	77.08%
Support Vector Machine (SVM)	78.84%	82.60%	73.63%
Random Forest	87.12%	83.06%	59.20%

http://fashion-mnist.s3-website.eu-central-l.amazonaws.com

Fine-Tuning the ResNet

- The best accuracy now is 93.42%
- Seems like transfer learning in our case is not that promising.

Models	Accuracy
Last layer fine tune	0.8617
Full fine tune	0.9342

Other Existing Models...

Fashion test accuracy	Preprocessing	Classifier
0.835	Crowd-sourced evaluation of human (with no fashion expertise) performance. 1000 randomly sampled test images, 3 labels per image, majority labelling.	Human Performance
0.926	HOG	HOG+SVM
0.898	scaling the pixel values to mean=0.0 and var=1.0	XgBoost
0.876	None	2 Conv Layers with max pooling (Keras)
0.916	None	2 Conv Layers with max pooling (Tensorflow) >300 epochs
0.903	None	2 Conv Layers with max pooling and ELU activation (PyTorch)
0.919	Normalization, random horizontal flip, random vertical flip, random translation, random rotation.	2 Conv Layers net
0.934	Augmentation, batch normalization	2 Conv Layers with 3 FC 500K parameters
0.926	None	3 Conv+pooling and 2 FC+dropout
0.936	Normalization and shift at most 2 pixel and horizontal flip	Capsule Network 8M parameters
0.947	standardization+augmentation+random erasing	CNN with optional shortcuts, dense-like connectivity
0.949	Normalization, random horizontal flip, random vertical flip, random translation, random rotation.	ResNet18
0.937	None	GoogleNet with cross-entropy loss
0.957	standard preprocessing (mean/std subtraction/division) and augmentation (random crops/horizontal flips)	Dual path network with wide resnet 28-10
0.935	None	VGG16 26M parameters
0.967	standard preprocessing (mean/std subtraction/division) and augmentation (random crops/horizontal flips)	WRN40-4 8.9M params

Q/A

Hong Kong University of Science and Technology Electronic & Computer Engineering Human Language Technology Center (HLTC)

Jason WU, Peng XU, Nayeon LEE