## Generative Models

#### Introduction

Yingzhen Li (yingzhen.li@imperial.ac.uk)

#### Supervised Learning

Data:  $(x_1, y_1), ..., (x_N, y_N) \sim p_{data}(x, y)$ 

Cat

Goal: learn a function to map  $x \rightarrow y$ 



Regression



Classification



GRASS, CAT, TREE, SKY Semantic Segmentation



**DOG**, **DOG**, **CAT** Object detection



# 95m

photos and videos are shared on Instagram

Instagram Business



of data will be created every day by 2025

IDC



tweets are sent every day

Twitter

# Most of the data are unlabelled

#### Unsupervised Learning

Data:  $x_1, ..., x_N \sim p_{data}(x)$  (no supervision signal)

Goal: inferring a function that describes the hidden structure of unlabelled data

Examples:

- Probability distribution/density estimation
- Dimensionality Reduction
- Clustering

All of them can be achieved by generative modelling!

#### Probability distribution/density estimation

Data:  $x_1, \dots, x_N \sim p_{data}(x)$ 

Goal: learn a distribution  $p_{\theta}(x) \approx p_{data}(x)$  with data  $x_1, \dots, x_N$ 



https://www.r-bloggers.com/2016/03/ballr-interactive-nba-shot-charts-with-r-and-shiny/

#### Generative Latent Variable Models

• Design  $p_{\theta}(x)$  as a generative latent variable model (LVM):

 $z \sim p_{\theta}(z), \qquad x \sim p_{\theta}(x|z)$  $\Rightarrow p_{\theta}(x) = \int p_{\theta}(x|z)p_{\theta}(z)dz$ 

*z*: latent variable (unobserved)

x: observation variable

*z*: digit label, writing style, ... *x*: hand-written digit



z: scene, viewing angle, lighting condition, ... x: photo image



Ζ

 $\boldsymbol{\chi}$ 

z: semantics, sentiment, ... x: generated text

• High-dimensional raw data are often sparse, perhaps lying on a low-dimensional manifold:



natural images vs all RGB images





• Principal Component Analysis (PCA):

Find principal components – orthogonal directions that capture most of the variance in the data

- 1<sup>st</sup> principal component direction of greatest variability
- 2<sup>nd</sup> principal component next orthogonal (uncorrelated) direction of greatest variability
- And so on ...

2<sup>nd</sup> PC 1<sup>st</sup> PC

Dimensionality reduction is achieved by projecting the data on the top K < d principal components ( $x \in R^d$ )

• Probabilistic Principal Component Analysis (Prob PCA):

 $p(z) = N(z; 0, I), z \in \mathbb{R}^{K}, K < d$  $p_{\theta}(x|z) = N(x; Wz, \sigma^{2}I), x \in \mathbb{R}^{d}$ 

- Parameters to optimize:  $\theta = W \in \mathbb{R}^{d \times K}$ , with the row vectors in W orthogonal to each other
- Trained using Maximum Likelihood
- Optimal *W* contains the top *K* principal components the top *K* eigenvectors of the data covariance matrix



Tipping and Bishop. Probabilistic principal component analysis. Journal of the Royal Statistical Society: Series B, 1999

- Auto-encoders for dimensionality reduction: •
  - Encoder network to extract data representations • (often with lower dimensionality)
  - Decoder network to reconstruct data given the representations ٠





#### networks trained by minimising reconstruction loss (e.g. L2 loss)



## Clustering

- Clustering: discover "group structure"
  - grouping datapoints into several clusters
  - Datapoints in the same cluster are similar
  - Datapoints in different clusters are "dissimilar"







gene data analysis

## Clustering

• Gaussian mixture model (GMM):

 $p_{\theta}(z) = Categorical(\pi),$  $\pi = (\pi_1, \dots, \pi_K), \pi_i = p_{\theta}(z = i), \sum_{i=1}^K \pi_i = 1$  $p_{\theta}(x|z) = N(x; \mu_z, \Sigma_z)$ 

- $z \in \{1, ..., K\}$ : index of the Gaussian component
- $\mu_z$ : mean of the *i*<sup>th</sup> Gaussian component if z = i
- $\Sigma_z$ : Covariance matrix of the  $i^{th}$  Gaussian component if z = i

 $\Rightarrow$  Clustering can be done by fitting a GMM model to the data





#### Representation learning

- Both dimensionality reduction and clustering can be viewed as representation learning
  - Hope: useful for downstream tasks

Representations used in downstream tasks:

• Classification (cat vs dog)



