

Unsupervised Learning

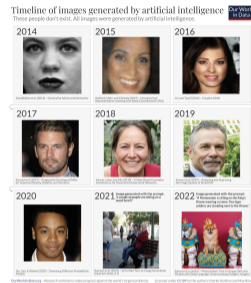
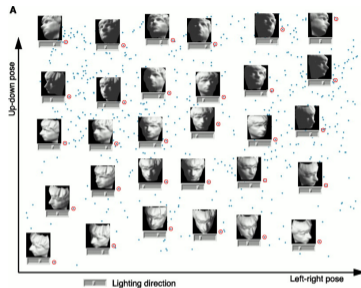
E. Le Pennec



Fall 2023

- 1 Unsupervised Learning?
- 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References

- 1 Unsupervised Learning?
- 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References



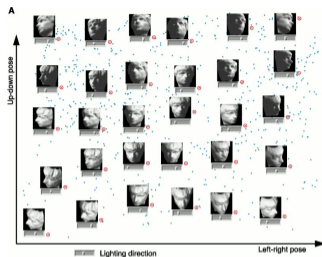
What is possible with data without labels?

- To group them?
- To visualize them in a 2 dimensional space?
- To generate more data?



To group them?

- **Data:** Base of customer data containing their properties and past buying records
- **Goal:** Use the customers *similarities* to find groups.
- **Clustering:** propose an explicit *grouping* of the customers
- **Visualization:** propose a representation of the customers so that the groups are *visible*. (Bonus)



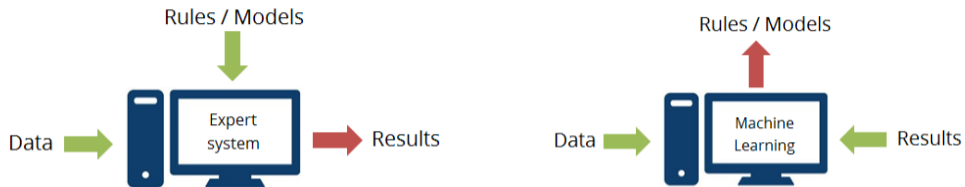
To visualize them?

- **Data:** Images of a single object
- **Goal:** Visualize the *similarities* between images.
- **Visualization:** propose a representation of the images so that similar images are *close*.
- **Clustering:** use this representation to cluster the images. (Bonus)



To generate more data?

- **Data:** Images.
- **Goal:** Generate images similar to the ones in the dataset.
- **Generative Modeling:** propose (and train) a generator.



The *classical* definition of Tom Mitchell

A computer program is said to learn from **experience E** with respect to some **class of tasks T** and **performance measure P**, if its performance at tasks in T, as measured by P, improves with experience E.

Experience, Task and Performance measure

- **Training data** : $\mathcal{D} = \{(\underline{X}_1, Y_1), \dots, (\underline{X}_n, Y_n)\}$ (i.i.d. $\sim \mathbb{P}$)
- **Predictor**: $f : \mathcal{X} \rightarrow \mathcal{Y}$ measurable
- **Cost/Loss function**: $\ell(f(\underline{X}), Y)$ measure how well $f(\underline{X})$ predicts Y
- **Risk**:

$$\mathcal{R}(f) = \mathbb{E}[\ell(Y, f(\underline{X}))] = \mathbb{E}_{\underline{X}} \left[\mathbb{E}_{Y|\underline{X}}[\ell(Y, f(\underline{X}))] \right]$$

- Often $\ell(f(\underline{X}), Y) = \|f(\underline{X}) - Y\|^2$ or $\ell(f(\underline{X}), Y) = \mathbf{1}_{Y \neq f(\underline{X})}$

Goal

- Learn a rule to construct a **predictor** $\hat{f} \in \mathcal{F}$ from the training data \mathcal{D}_n s.t. **the risk** $\mathcal{R}(\hat{f})$ is **small on average** or with high probability with respect to \mathcal{D}_n .

Experience, Task and Performance measure

- **Training data** : $\mathcal{D} = \{\underline{X}_1, \dots, \underline{X}_n\}$ (i.i.d. $\sim \mathbb{P}$)
 - **Task**: ???
 - **Performance measure**: ???
- No obvious task definition!

Tasks for this lecture

- **Dimension reduction**: construct a map of the data in a **low dimensional** space without **distorting** it too much.
- **Clustering (or unsupervised classification)**: construct a **grouping** of the data in **homogeneous** classes.
- **Generative modeling**: **generate** new conditional samples.

- **Training data** : $\mathcal{D} = \{\underline{X}_1, \dots, \underline{X}_n\} \in \mathcal{X}^n$ (i.i.d. $\sim \mathbb{P}$)
- Space \mathcal{X} of possibly high dimension.

Dimension Reduction Map

- Construct a map Φ from the space \mathcal{X} into a space \mathcal{X}' of **smaller dimension**:

$$\Phi : \mathcal{X} \rightarrow \mathcal{X}'$$

$$\underline{X} \mapsto \Phi(\underline{X})$$

- Map can be defined only on the dataset.

Motivations

- Visualization of the data
- Dimension reduction (or embedding) before further processing

- Need to control the **distortion** between \mathcal{D} and $\Phi(\mathcal{D}) = \{\Phi(\underline{X}_1), \dots, \Phi(\underline{X}_n)\}$

Distortion(s)

- Reconstruction error:
 - Construct $\tilde{\Phi}$ from \mathcal{X}' to \mathcal{X}
 - Control the error between \underline{X} and its reconstruction $\tilde{\Phi}(\Phi(\underline{X}))$
 - Relationship preservation:
 - Compute a *relation* \underline{X}_i and \underline{X}_j and a *relation* between $\Phi(\underline{X}_i)$ and $\Phi(\underline{X}_j)$
 - Control the difference between those two *relations*.
-
- Lead to different constructions. . . .

- **Training data** : $\mathcal{D} = \{\underline{X}_1, \dots, \underline{X}_n\} \in \mathcal{X}^n$ (i.i.d. $\sim \mathbb{P}$)
- Latent groups?

Clustering

- Construct a map f from \mathcal{D} to $\{1, \dots, K\}$ where K is a number of classes to be fixed:

$$f : \underline{X}_i \mapsto k_i$$

- Similar to classification except:
 - no ground truth (no given labels)
 - label only elements of the dataset!

Motivations

- Interpretation of the groups
- Use of the groups in further processing

- Need to define the **quality** of the cluster.
- No obvious measure!

Clustering quality

- Inner homogeneity: samples in the same group should be similar.
- Outer inhomogeneity: samples in two different groups should be different.
- Several possible definitions of similar and different.
- Often based on the distance between the samples.
- Example based on the Euclidean distance:
 - Inner homogeneity = intra-class variance,
 - Outer inhomogeneity = inter-class variance.
- **Beware:** choice of the number of clusters K often complex!

- **Training data** : $\mathcal{D} = \{(\underline{X}_1, \underline{Y}_1), \dots, (\underline{X}_n, \underline{Y}_n)\} \in (\mathcal{X} \times \mathcal{Y})^n$ (i.i.d. $\sim \mathbb{P}$)
- Same kind of data than for supervised learning if $\mathcal{Y} \neq \emptyset$.

Generative Modeling

- Construct a map G from the product of \mathcal{Y} and a randomness source Ω to \mathcal{X}

$$G : \mathcal{Y} \times \Omega \rightarrow \mathcal{X}$$

$$(Y, \omega) \mapsto X$$

- Unconditional model if $\mathcal{Y} = \emptyset$...

Motivation

- Generate plausible novel conditional samples based on a given dataset.

Sample Quality

- Related to the proximity between the law of $G(Y, \omega)$ and the law of $X|Y$.
- Most classical choice is the Kullback-Leibler divergence.

Ingredients

- Generator $G_\theta(Y, \omega)$ and cond. density prob. $P_\theta(X|Y)$ (Explicit vs implicit link)
- Simple / Complex / Approximate estimation...

Some Possible Choices

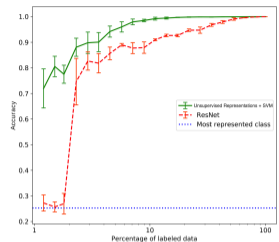
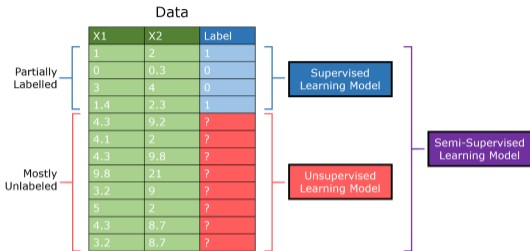
	Probabilistic model	Generator	Estimation
Base	Simple (parametric)	Explicit	Simple (ML)
Flow	Image of simple model	Explicit	Simple (ML)
Factorization	Factorization of simple model	Explicit	Simple (ML)
VAE	Simple model with latent var.	Explicit	Approximate (ML)
EBM	Arbitrary	Implicit (MCMC)	Complex (ML/score/discrim.)
Diffusion	Continuous noise	Implicit (MCMC)	Complex (score)
	Discrete Noise with latent var.	Explicit	Approximate (ML)
GAN	Implicit	Explicit	Complex (Discrimination)

- SOTA: Diffusion based approach!

- **General observation:** most data do not have a label !
- **Example:** The number of images on which someone has described the content of the image is a *tiny fraction* of the images online.
- Labeling is very expensive and time consuming
- A lot of information can be extracted from the structure of the data, before seeing any label.

How can we leverage the large quantity of un-labeled data?

- Learn relevant features (= **representations**) in an unsupervised fashion
 - Use those features to solve a supervised task with a fraction of labeled data.
-
- **Semi-supervised framework**
 - ↗ Very useful in practice, for images, time series, text.



Semu-Supervised Framework

- With representation learned in an unsupervised fashion + a simple linear model, one can achieve the same performance with 10% of data labeled than with a fully annotated dataset.
- Complementary regularization based approaches also exist.

Except for generative modeling, the learner is always right

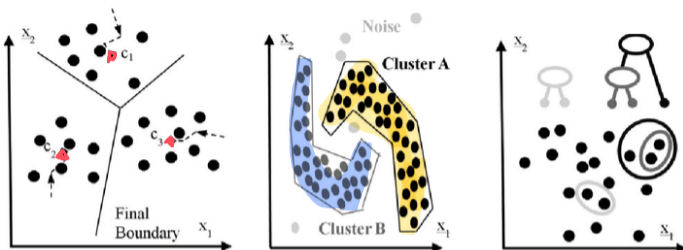
- A subjective measure of performance
 - Subjective choices for the algorithmic constraints (e.g., the type of transformation of the data we allow for low-dimensional representation, type of groups in clustering)
 - \Rightarrow Very difficult or impossible to tell which is the *best* method.
-
- Yet:
 - Extremely important in practice:
 - 90-99% of the data is un-labeled!
 - the tasks themselves are fundamental
 - Huge success in various fields (NLP, images. . .)

Lecture goals for the three main tasks

- Discussing possible choices of measures of performance and algorithmic constraints
- Understand the correspondences between those choices and a variety of classical algorithms
- For the simplest algorithms (PCA, k-means), get a precise mathematical understanding of the learning process.

- 1 Unsupervised Learning?
- 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References

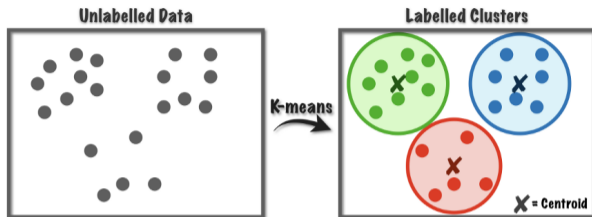
- 1 Unsupervised Learning?
- 2 **A First Glimpse**
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References



- No simple or unanimous definition!
- Require a notion of similarity/difference. . .

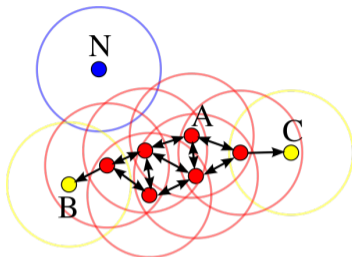
Three main approaches

- A group is a set of samples similar to a prototype.
- A group is a set of samples that can be linked by contiguity.
- A group can be obtained by fusing some smaller groups. . .



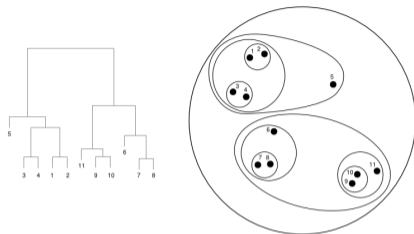
Prototype Approach

- A group is a set of samples similar to a prototype.
- Most classical instance: *k*-means algorithm.
- Principle: alternate prototype choice for the current groups and group update based on those prototypes.
- Number of groups fixed at the beginning
- No need to compare the samples between them!



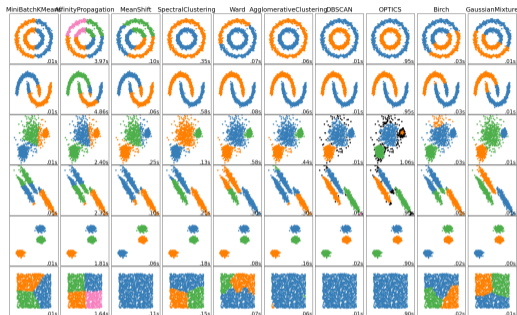
Contiguity Approach

- A group is the set of samples that can be linked by contiguity.
- Most classical instance: DBScan
- Principle: group samples by contiguity if possible (proximity and density)
- Some samples may remain isolated.
- Number of groups controlled by the scale parameter.



Agglomerative Approach

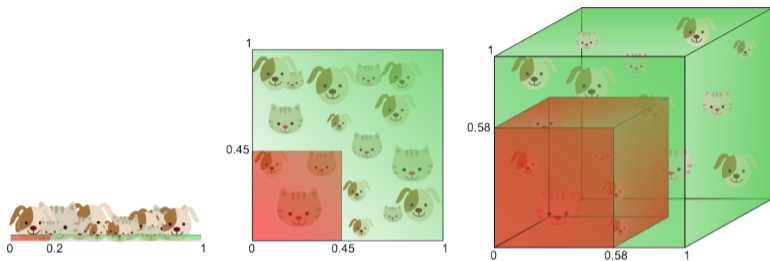
- A group can be obtained by fusing some smaller groups. . .
- Hierarchical clustering principle: sequential merging of groups according to a *best merge* criterion
- Numerous variations on the merging criterion. . .
- Number of groups chosen afterward.



No method or number of groups is better than the others...

- Criterion not necessarily explicit!
- No cross validation possible
- Choice of the number of groups: a priori, heuristic, *based on the final usage...*

- 1 Unsupervised Learning?
- 2 **A First Glimpse**
 - Clustering
 - **Dimensionality Curse**
 - Dimension Reduction
 - Generative Modeling
- 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References



- **DISCLAIMER: Even if they are used everywhere beware of the usual distances in high dimension!**

Dimensionality Curse

- Previous approaches based on distances.
- Surprising behavior in high dimension: everything is ((often) as) far away.
- Beware of categories. . .

- **DISCLAIMER: Even if they are used everywhere beware of the usual distances in high dimension!**

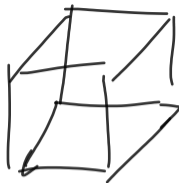
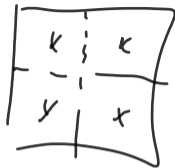
High Dimensional Geometry Course

- Folks theorem: In high dimension, everyone is alone.
- Theorem: If $\underline{X}_1, \dots, \underline{X}_n$ in the hypercube of dimension d such that their coordinates are i.i.d then

$$d^{-1/p} \left(\max \|\underline{X}_i - \underline{X}_j\|_p - \min \|\underline{X}_i - \underline{X}_j\|_p \right) = 0 + O_P \left(\sqrt{\frac{\log n}{d}} \right)$$
$$\frac{\min \|\underline{X}_i - \underline{X}_j\|_p}{\max \|\underline{X}_i - \underline{X}_j\|_p} = 1 + O_P \left(\sqrt{\frac{\log n}{d}} \right).$$

- When d is large, all the points are almost equidistant. . .
- Nearest neighbors are meaningless!

$$n \leq 2^d$$



d 1

2

3

d

m 2

4

8

2^d

- 1 Unsupervised Learning?
- 2 **A First Glimpse**
 - Clustering
 - Dimensionality Curse
 - **Dimension Reduction**
 - Generative Modeling
- 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References

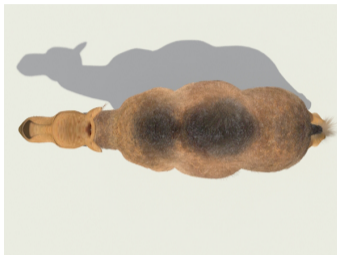
Visualization and Dimension Reduction

- How to view a dataset in high dimension !
- High dimension: dimension larger than 2!
- **Projection** onto a 2D space.



Visualization and Dimension Reduction

- How to view a dataset in high dimension !
- High dimension: dimension larger than 2!
- **Projection** onto a 2D space.



Visualization and Dimension Reduction

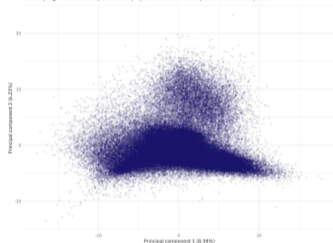
- How to view a dataset in high dimension !
- High dimension: dimension larger than 2!
- **Projection** onto a 2D space.



Visualization and Dimension Reduction

- How to view a dataset in high dimension !
- High dimension: dimension larger than 2!
- **Projection** onto a 2D space.

Projection of Stack Overflow traffic on to the first two principal components
The very high dimensional space can be projected down onto components we have explored

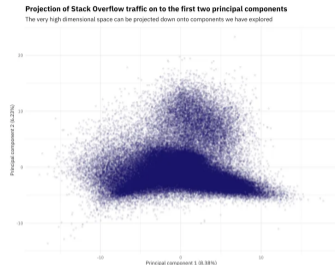


proportion
orthogonale

- Simple formula: $\tilde{X} = P(X - m)$

How to chose P ?

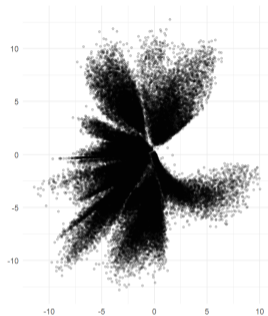
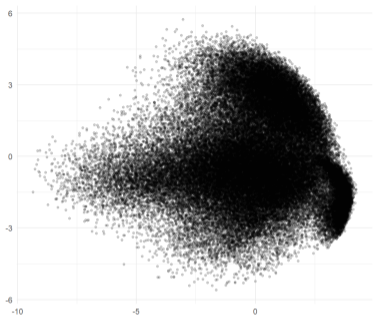
- Maximising the dispersion of the points?
- Allowing to well reconstruct X from \tilde{X} ?
- Preserving the relationship between the X through those between the \tilde{X} ?



- Simple formula: $\tilde{X} = P(X - m)$

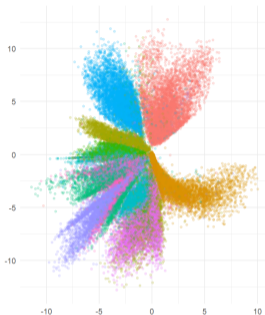
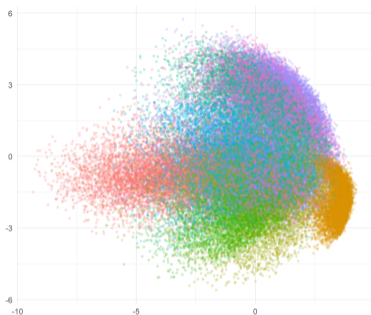
How to chose P ?

- Maximising the dispersion of the points?
- Allowing to well reconstruct X from \tilde{X} ?
- Preserving the relationship between the X through those between the \tilde{X} ?
- The 3 approaches yield the same solution!



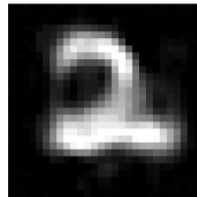
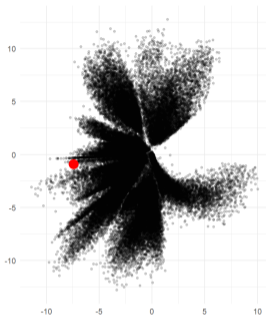
Reconstruction Approaches

- Learn a formula to encode and one formula to decode.
- Auto-encoder structure
- Yields a formula for new points.



Reconstruction Approaches

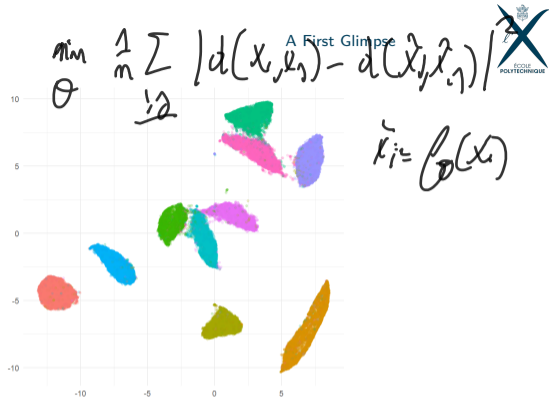
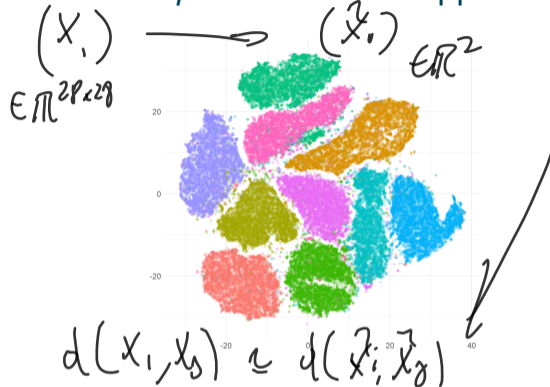
- Learn a formula to encode and one formula to decode.
- Auto-encoder structure
- Yields a formula for new points.



Reconstruction Approaches

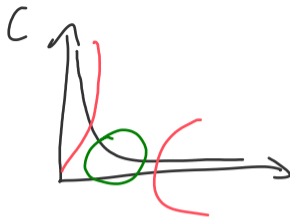
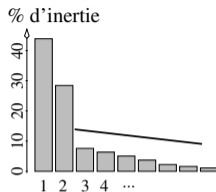
- Learn a formula to encode and one formula to decode.
- Auto-encoder structure
- Yields a formula for new points.

Relationship Preservation Approaches



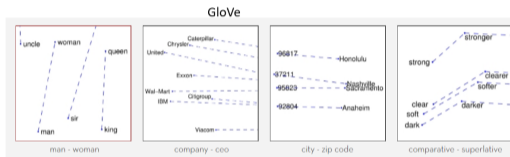
Relationship Preservation Approaches

- Based on the definition of the relationship notion (in both worlds).
- Huge flexibility
- Not always yields a formula for new points.



No Better Choice?

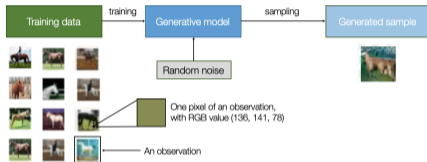
- Different criterion for different methods: impossible to use cross-validation.
 - The larger the dimension the easier is to be faithful!
 - In visualization, dimension 2 is the only choice.
 - Heuristic criterion for the dimension choice: elbow criterion (no more gain), stability...
-
- Dimension Reduction is rarely used standalone but rather as a step in a predictive/prescriptive method.
 - The dimension becomes an hyper-parameter of this method.



Representation Learning

- How to transform arbitrary objects into numerical vectors?
- Objects: Categorical variables, Words, Images/Sounds. . .
- The two previous dimension reduction approaches can be used (given possibly a first simple high dimensional representation)

- 1 Unsupervised Learning?
- 2 **A First Glimpse**
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - **Generative Modeling**
- 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References

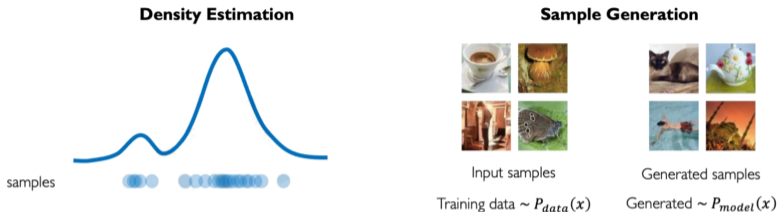


Timeline of images generated by artificial intelligence. These people don't exist. All images were generated by artificial intelligence.



Generative Modeling

- Generate new samples similar to the ones in an original dataset.
- Generation may be conditioned by an input.
- Key for image generation. . . and chatbot!

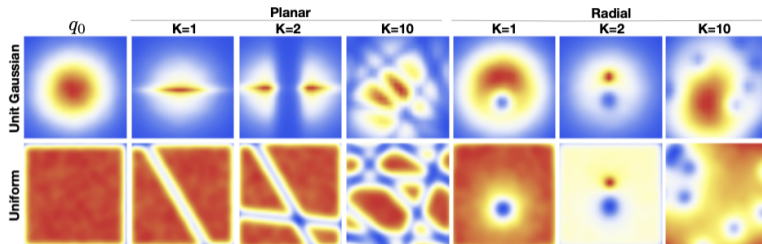


How can we learn $P_{model}(x)$ similar to $P_{data}(x)$?

- **Heuristic:** If we can estimate the (conditional) law \mathbb{P} of the data and can simulate it, we can obtain new samples similar to the input ones.

Estimation and Simulation

- How to estimate the density?
- How to simulate the estimate density?
- Other possibilities?



Parametric Model, Image and Factorization

- Use

- a simple parametric model,...
- or the image of a parametric model (flow),...
- or a factorization of a parametric model (recurrent model)

as they are *simple* to estimate and to simulate.

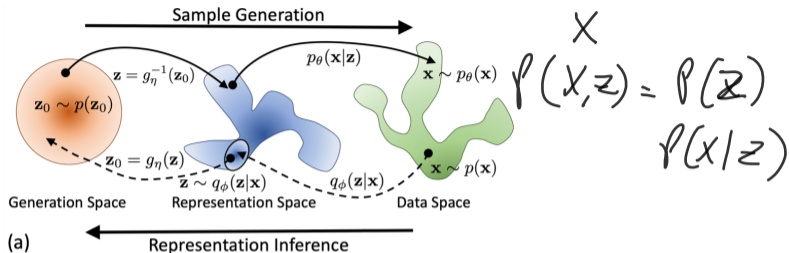
$$p_\theta$$

$$p_\theta(p_\theta)$$

$$f(x)$$

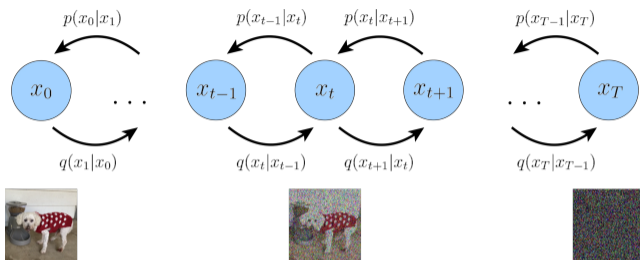
$$\prod P(x_i | x_{i-1}, \dots, x_1)$$

- Estimation by Maximum Likelihood principle.
- Recurrent models are used in Large Language Models!



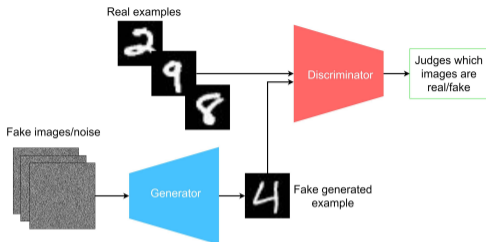
Latent Variable

- Generate first a (low dimensional) latent variable Z from which X is easy to sample.
- Estimation based on approximate Maximum Likelihood (VAE/ELBO)
- The latent variable can be generated by a simple method (or a more complex one...).



Monte Carlo Markov Chain

- Rely on much more complex probability model. . .
- which can only be simulated numerically.
- Often combined with noise injection to stabilize the numerical scheme (Diffusion).
- Much more expensive to simulate than with Latent Variable approaches.



Generative Adversarial Network

- Bypass the density estimation problem, by transforming the problem into a competition between the generator and a discriminator.
- The better the generator, the harder it is for the generator to distinguish true samples from synthetic ones.
- No explicit density!
- Fast simulator but **unstable training...**

- 1 Unsupervised Learning?
- 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 3 **Dimension Reduction**
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References

- **Training data** : $\mathcal{D} = \{\underline{X}_1, \dots, \underline{X}_n\} \in \mathcal{X}^n$ (i.i.d. $\sim \mathbb{P}$)
- Space \mathcal{X} of possibly high dimension.

Dimension Reduction Map

- Construct a map Φ from the space \mathcal{X} into a space \mathcal{X}' of **smaller dimension**:

$$\Phi : \mathcal{X} \rightarrow \mathcal{X}'$$

$$\underline{X} \mapsto \Phi(\underline{X})$$

Criterion

- Reconstruction error
- Relationship preservation

- 1 Unsupervised Learning?
- 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 3 **Dimension Reduction**
 - **Simplification**
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References

A Projection Based Approach

- Observations: $\underline{X}_1, \dots, \underline{X}_n \in \mathbf{R}^d$
- Simplified version: $\Phi(\underline{X}_1), \dots, \Phi(\underline{X}_n) \in \mathbf{R}^d$ with Φ an affine projection preserving the mean $\Phi(\underline{X}) = P(\underline{X} - m) + m$ with $P^\top = P = P^2$ and $m = \frac{1}{n} \sum_i \underline{X}_i$.

How to choose P ?

- **Inertia criterion:**

$$\max_P \sum_{i,j} \|\Phi(\underline{X}_i) - \Phi(\underline{X}_j)\|^2?$$

- **Reconstruction criterion:**

$$\min_P \sum_i \|\underline{X}_i - \Phi(\underline{X}_i)\|^2?$$

- **Relationship criterion:**

$$\min_P \sum_{i,j} |(\underline{X}_i - m)^\top (\underline{X}_j - m) - (\Phi(\underline{X}_i) - m)^\top (\Phi(\underline{X}_j) - m)|^2?$$

- **Rk:** Best solution is $P = I$! Need to reduce the rank of the projection to $d' < d \dots$

- **Heuristic:** a good representation is such that the projected points are far apart.

Two views on inertia

- Inertia:

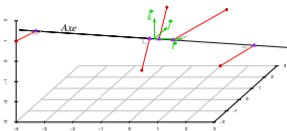
$$I = \frac{1}{2n^2} \sum_{i,j} \|\underline{X}_i - \underline{X}_j\|^2 = \frac{1}{n} \sum_{i=1}^n \|\underline{X}_i - m\|^2$$

- 2 times the mean squared distance to the mean = Mean squared distance between individual

Inertia criterion (Principal Component Analysis)

- Criterion: $\max_P \sum_{i,j} \frac{1}{2n^2} \|P\underline{X}_i - P\underline{X}_j\|^2 = \max_P \frac{1}{n} \sum_i \|P\underline{X}_i - m\|^2$

- **Solution:** Choose P as a projection matrix on the space spanned by the d' first eigenvectors of $\Sigma = \frac{1}{n} \sum_i (\underline{X}_i - m)(\underline{X}_i - m)^\top$



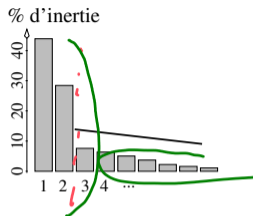
- $\tilde{X} = m + a^\top (X - m)a$ with $\|a\| = 1$
- Inertia: $\frac{1}{n} \sum_{i=1}^n a^\top (X_i - m)(X_i - m)^\top a$

Principal Component Analysis: optimization of the projection

- Maximization of $\tilde{l} = \frac{1}{n} \sum_{i=1}^n a^\top (X_i - m)(X_i - m)^\top a = a^\top \Sigma a$ with

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (X_i - m)(X_i - m)^\top \text{ the empirical covariance matrix.}$$

- Explicit optimal choice given by the eigenvector of the largest eigenvalue of Σ .



Principal Component Analysis : sequential optimization of the projection

- Explicit optimal solution obtain by the projection on the eigenvectors of the largest eigenvalues of Σ .
- Projected inertia given by the sum of those eigenvalues.
- Often fast decay of the eigenvalues: some dimensions are much more important than others.
- Not exactly the curse of dimensionality setting. . .
- Yet a lot of *small* dimension can drive the distance!

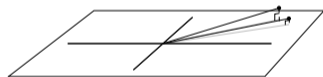
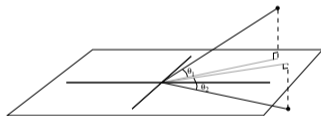
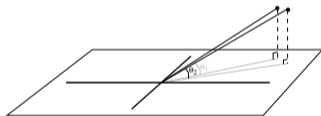
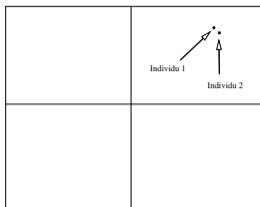
- **Heuristic:** a good representation is such that the projected points are close to the original ones.

Reconstruction Criterion

- Criterion: $\min_P \sum_i \frac{1}{n} \|\underline{X}_i - (P(\underline{X}_i - m) + m)\|^2 = \min_P \frac{1}{n} \sum_i \|(I - P)(\underline{X}_i - m)\|^2$
- **Solution:** Choose P as a projection matrix on the space spanned by the d' first eigenvectors of $\Sigma = \frac{1}{n} \sum_i (\underline{X}_i - m)(\underline{X}_i - m)^\top$

- Same solution with a different heuristic!
- Proof (Pythagora):

$$\sum_i \|\underline{X}_i - m\|^2 = \sum_i \left(\|P(\underline{X}_i - m)\|^2 + \|(I - P)(\underline{X}_i - m)\|^2 \right)$$



Close projection doesn't mean close individuals!

- Same projections but different situations.
- Quality of the reconstruction measured by the angle with the projection space!

- **Heuristic:** a good representation is such that the projected points scalar products are similar to the original ones.

Relationship Criterion (Multi Dimensional Scaling)

- Criterion: $\min_P \sum_{i,j} |(\underline{X}_i - m)^\top (\underline{X}_j - m) - (\Phi(\underline{X}_i) - m)^\top (\Phi(\underline{X}_j) - m)|^2$
- **Solution:** Choose P as a projection matrix on the space spanned by the d' first eigenvectors of $\Sigma = \frac{1}{n} \sum_i (\underline{X}_i - m)(\underline{X}_i - m)^\top$
- Same solution with a different heuristic!
- Much more involved justification!

- PCA model: $\underline{X} - m \simeq P(\underline{X} - m)$
- **Prop:** $P = VV^T$ with V an orthormal family in dimension d of size d' .
- PCA model with V : $\underline{X} - m \simeq VV^T(\underline{X} - m)$ where $\tilde{\underline{X}} = V^T(\underline{X} - m) \in \mathbb{R}^{d'}$
- Row vector rewriting: $\underline{X}^T - m^T \simeq \tilde{\underline{X}}^T V^T$

Matrix Rewriting and Low Rank Factorization

- Matrix rewriting

$$\begin{array}{ccc}
 \begin{array}{|c} \hline \underline{X}_1^T - m^T \\ \hline \vdots \\ \hline \underline{X}_n^T - m^T \\ \hline \end{array} & \simeq & \begin{array}{|c} \hline \tilde{\underline{X}}_1^T \\ \hline \vdots \\ \hline \tilde{\underline{X}}_n^T \\ \hline \end{array} \begin{array}{|c} \hline \mathbf{V}^T \\ \hline \end{array} \\
 (n \times d) & & (n \times d') \quad (d' \times d)
 \end{array}$$

- Low rank matrix factorization! (Truncated SVD solution...)

SVD Decomposition

- Any matrix $n \times d$ matrix A can be decomposed as

$$\begin{array}{c} \boxed{\mathbf{A}} \\ (n \times d) \end{array} = \begin{array}{c} \boxed{\mathbf{U}} \\ (n \times n) \end{array} \begin{array}{c} \boxed{D} \\ (n \times d) \end{array} \begin{array}{c} \boxed{\mathbf{W}^T} \\ (d \times d) \end{array}$$

with U and W two orthonormal matrices and D a *diagonal* matrix with decreasing values.

Low Rank Approximation

- The best low rank approximation or rank r is obtained by restriction of the matrices to the first r dimensions:

$$\begin{array}{c}
 \boxed{\mathbf{A}} \\
 (n \times d)
 \end{array}
 \simeq
 \begin{array}{c}
 \boxed{\mathbf{U}_r} \\
 (n \times r)
 \end{array}
 \begin{array}{c}
 \boxed{D_{r,r}} \\
 (r \times r)
 \end{array}
 \begin{array}{c}
 \boxed{\mathbf{W}_r^\top} \\
 (r \times d)
 \end{array}$$

for both the operator norm and the Frobenius norm!

- PCA: Low rank approximation with Frobenius norm, $d' = r$ and

$$\begin{pmatrix} \underline{X}_1^\top - m^\top \\ \vdots \\ \underline{X}_n^\top - m^\top \end{pmatrix} \leftrightarrow A, \quad \begin{pmatrix} \tilde{X}_1^\top \\ \vdots \\ \tilde{X}_n^\top \end{pmatrix} \leftrightarrow \mathbf{U}_r D_{r,r}, \quad \mathbf{V}^\top \leftrightarrow \mathbf{W}_r^\top$$

SVD Decompositions

- Recentered data:

$$\mathbf{R} = \begin{pmatrix} \underline{X}_1^\top - m^\top \\ \vdots \\ \underline{X}_n^\top - m^\top \end{pmatrix} = \mathbf{U}\mathbf{D}\mathbf{W}^\top$$

- Covariance matrix:

$$\Sigma = \mathbf{R}^\top \mathbf{R} = \mathbf{W}\mathbf{D}^\top \mathbf{D}\mathbf{W}$$

$$\Sigma_{\ell\ell} = \frac{1}{n} \sum_i X_i^{(a)} X_i^{(e)}$$

with $\mathbf{D}^\top \mathbf{D}$ diagonal.

- Gram matrix (matrix of scalar products):

$$\mathbf{G} = \mathbf{R}\mathbf{R}^\top = \mathbf{U}\mathbf{D}\mathbf{D}^\top \mathbf{U}$$

$$G_{ij} = \frac{1}{\epsilon} \sum_{\ell} X_i^{(a)} X_j^{(a)}$$

with $\mathbf{D}\mathbf{D}^\top$ diagonal.

- Those are the same \mathbf{U} , \mathbf{W} and \mathbf{D} , hence the link between all the approaches.

- 1 Unsupervised Learning?
- 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 3 **Dimension Reduction**
 - Simplification
 - **Reconstruction Error**
 - Relationship Preservation
 - Comparing Methods?
- 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References

Goal

- Construct a map Φ from the space \mathcal{X} into a space \mathcal{X}' of **smaller dimension**:

$$\Phi : \mathcal{X} \rightarrow \mathcal{X}'$$

$$\underline{X} \mapsto \Phi(\underline{X})$$

- Construct $\tilde{\Phi}$ from \mathcal{X}' to \mathcal{X}

- Control the error between \underline{X} and its reconstruction $\tilde{\Phi}(\Phi(\underline{X}))$

- Canonical example for $\underline{X} \in \mathbb{R}^d$: find Φ and $\tilde{\Phi}$ in a parametric family that minimize

$$\frac{1}{n} \sum_{i=1}^n \|\underline{X}_i - \tilde{\Phi}(\Phi(\underline{X}_i))\|^2$$

- $\mathcal{X} \in \mathbb{R}^d$ and $\mathcal{X}' = \mathbb{R}^{d'}$
- Affine model $\underline{X} \sim m + \sum_{l=1}^{d'} \underline{X}'^{(l)} V^{(l)}$ with $(V^{(l)})$ an orthonormal family.
- Equivalent to:

$$\Phi(\underline{X}) = V^\top (\underline{X} - m) \quad \text{and} \quad \tilde{\Phi}(\underline{X}') = m + V \underline{X}'$$

- Reconstruction error criterion:

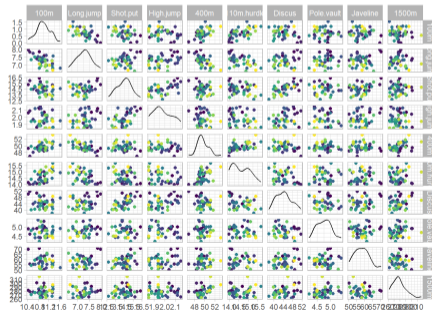
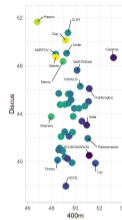
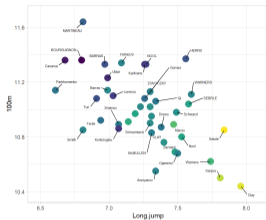
$$\frac{1}{n} \sum_{i=1}^n \|\underline{X}_i - (m + VV^\top (\underline{X}_i - m))\|^2$$

- **Explicit solution:** m is the empirical mean and V is any orthonormal basis of the space spanned by the d' first eigenvectors (the one with largest eigenvalues) of the empirical covariance matrix $\frac{1}{n} \sum_{i=1}^n (\underline{X}_i - m)(\underline{X}_i - m)^\top$.

PCA Algorithm

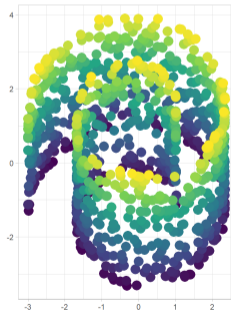
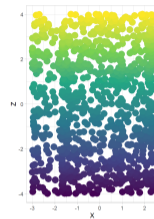
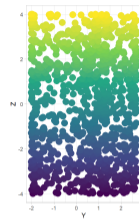
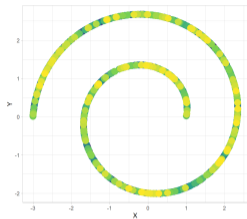
- Compute the empirical mean $m = \frac{1}{n} \sum_{i=1}^n \underline{X}_i$
 - Compute the empirical covariance matrix $\frac{1}{n} \sum_{i=1}^n (\underline{X}_i - m)(\underline{X}_i - m)^\top$.
 - Compute the d' first eigenvectors of this matrix: $V^{(1)}, \dots, V^{(d')}$
 - Set $\Phi(\underline{X}) = V^\top (\underline{X} - m)$
-
- Complexity: $O(n(d + d^2) + d'd^2)$
 - Interpretation:
 - $\Phi(\underline{X}) = V^\top (\underline{X} - m)$: coordinates in the restricted space.
 - $V^{(i)}$: influence of each original coordinates in the i th new one.
 - **Scaling:** This method is not invariant to a scaling of the variables! It is custom to normalize the variables (at least within groups) before applying PCA.

Decathlon

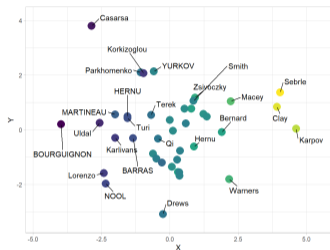


Swiss Roll

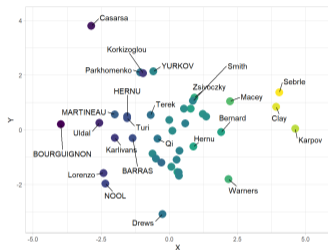
Dimension Reduction



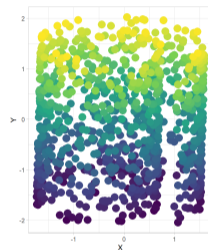
Principal Component Analysis



Decathlon



Decathlon
Renormalized



Swiss Roll

△ Normalisation
automatique

Multiple Factor Analysis

Dimension Reduction



- PCA assumes $\mathcal{X} = \mathbb{R}^d$
- How to deal with categorical values?
- MFA = PCA with clever coding strategy for categorical values.

$$\{1, \dots, V\} \rightarrow \left\{ \begin{array}{l} (1, 0, 0, \dots, 0) \\ (0, 1, 0, \dots, 0) \\ \vdots \\ (0, \dots, 0, 1) \end{array} \right\}$$

Categorical value code for a single variable

- Classical redundant dummy coding:

$$\underline{X} \in \{1, \dots, V\} \mapsto P(\underline{X}) = (\mathbf{1}_{\underline{X}=1}, \dots, \mathbf{1}_{\underline{X}=V})^\top$$

- Compute the mean (i.e. the empirical proportions): $\bar{P} = \frac{1}{n} \sum_{i=1}^n P(\underline{X}_i)$

- Renormalize $P(\underline{X})$ by $1/\sqrt{(V-1)\bar{P}}$:

$$P(\underline{X}) = (\mathbf{1}_{\underline{X}=1}, \dots, \mathbf{1}_{\underline{X}=V}) \mapsto \left(\frac{\mathbf{1}_{\underline{X}=1}}{\sqrt{(V-1)\bar{P}_1}}, \dots, \frac{\mathbf{1}_{\underline{X}=V}}{\sqrt{(V-1)\bar{P}_V}} = P^r(\underline{X}) \right)$$

- χ^2 type distance!

- PCA becomes the minimization of

$$\frac{1}{n} \sum_{i=1}^n \|P^r(\underline{X}_i) - (m + VV^\top (P^r(\underline{X}_i) - m))\|^2$$

$$= \frac{1}{n} \sum_{i=1}^n \sum_{v=1}^V \frac{|\mathbf{1}_{\underline{X}_i=v} - (m' + \sum_{l=1}^{d'} V^{(l)\top} (P(\underline{X}_i) - m') V^{(l,v)})|^2}{(V-1)\bar{P}_v}$$

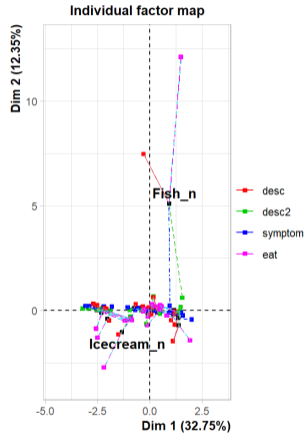
- Interpretation:

- $m' = \bar{P}$
- $\Phi(\underline{X}) = V^\top (P^r(\underline{X}) - m)$: coordinates in the restricted space.
- $V^{(l)}$ can be interpreted as a probability profile.

- Complexity: $O(n(V + V^2) + d'V^2)$
- Link with Correspondence Analysis (CA)

MFA Algorithm

- Redundant dummy coding of each categorical variable.
 - Renormalization of each block of dummy variable.
 - Classical PCA algorithm on the resulting variables
-
- Interpretation as a reconstruction error with a rescaled/ χ^2 metric.
 - Interpretation:
 - $\Phi(\underline{X}) = V^\top (P^r(\underline{X}) - m)$: coordinates in the restricted space.
 - $V^{(l)}$: influence of each modality/variable in the l th new coordinates.
 - **Scaling:** This method is not invariant to a scaling of the continuous variables! It is custom to normalize the variables (at least within groups) before applying PCA.



PCA Model

- PCA: Linear model assumption

$$\underline{X} \simeq m + \sum_{l=1}^{d'} \underline{X}'^{(l)} V^{(l)} = m + V \underline{X}'$$

- with

- $V^{(l)}$ orthonormal
- $\underline{X}'^{(l)}$ without constraints.

- Two directions of extension:
 - Other constraints on V (or the coordinates in the restricted space): ICA, NMF, Dictionary approach
 - PCA on a non-linear image of \underline{X} : kernel-PCA
- Much more complex algorithm!

ICA (Independent Component Analysis)

- Linear model assumption

$$\underline{X} \simeq m + \sum_{l=1}^{d'} \underline{X}^{',(l)} V^{(l)} = m + V \underline{X}'$$

- with

- $V^{(l)}$ without constraints.
- $\underline{X}^{',(l)}$ *independent*

NMF (Non Negative Matrix Factorization)

- (Linear) Model assumption

$$\underline{X} \simeq \sum_{l=1}^{d'} \underline{X}^{',(l)} V^{(l)} = V \underline{X}'$$

- with

- $V^{(l)}$ non-negative
- $\underline{X}^{',(l)}$ non-negative.

Dictionary

- (Linear) Model assumption

$$\underline{X} \simeq m + \sum_{l=1}^{d'} \underline{X}'^{(l)} V^{(l)} = m + V \underline{X}'$$

- with
 - $V^{(l)}$ without constraints
 - \underline{X}' sparse (with a lot of 0)

kernel PCA

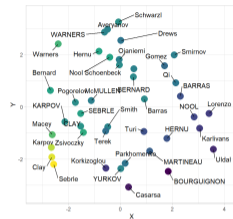
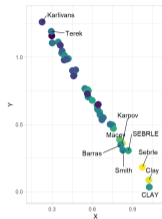
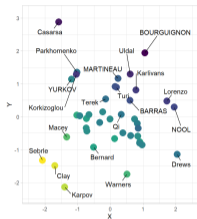
- Linear model assumption

$$\Psi(\underline{X} - m) \simeq \sum_{l=1}^{d'} \underline{X}'^{(l)} V^{(l)} = V \underline{X}'$$

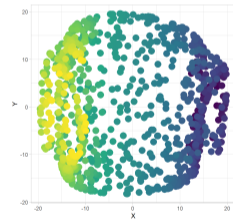
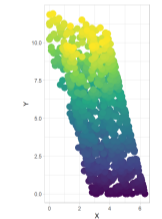
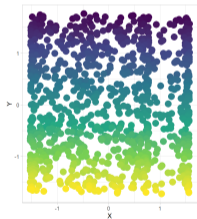
- with
 - $V^{(l)}$ orthonormal
 - \underline{X}'_l without constraints.

Non Linear PCA

Decathlon



Swiss Roll



ICA

NMF

Kernel PCA

Deep Auto Encoder

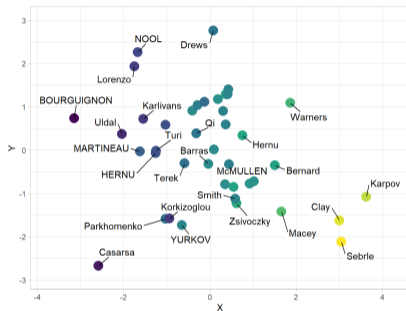
- Construct a map Φ with a **NN** from the space \mathcal{X} into a space \mathcal{X}' of smaller dimension:

$$\begin{aligned}\Phi : \mathcal{X} &\rightarrow \mathcal{X}' \\ \underline{X} &\mapsto \Phi(\underline{X})\end{aligned}$$

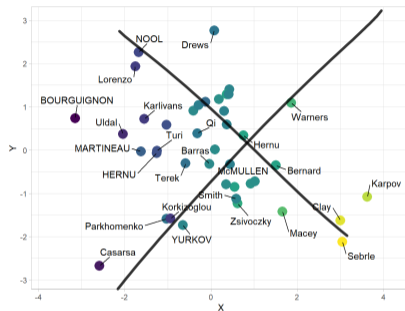
- Construct $\tilde{\Phi}$ with a **NN** from \mathcal{X}' to \mathcal{X}
- Control the error between \underline{X} and its reconstruction $\tilde{\Phi}(\Phi(\underline{X}))$:

$$\frac{1}{n} \sum_{i=1}^n \|\underline{X}_i - \tilde{\Phi}(\Phi(\underline{X}_i))\|^2$$

- Optimization by gradient descent.
- NN can be replaced by another parametric function...



Shallow Auto Encoder
(PCA)



Deep Auto Encoder

- 1 Unsupervised Learning?
- 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 3 **Dimension Reduction**
 - Simplification
 - Reconstruction Error
 - **Relationship Preservation**
 - Comparing Methods?
- 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References

- Different point of view!
- Focus on pairwise relation $\mathcal{R}(\underline{X}_i, \underline{X}_j)$.

Distance Preservation

- Construct a map Φ from the space \mathcal{X} into a space \mathcal{X}' of **smaller dimension**:

$$\Phi: \mathcal{X} \rightarrow \mathcal{X}'$$

$$\underline{X} \mapsto \Phi(\underline{X}) = \underline{X}'$$

- such that

$$\mathcal{R}(\underline{X}_i, \underline{X}_j) \sim \mathcal{R}'(\underline{X}'_i, \underline{X}'_j)$$

- Most classical version (MDS):

- Scalar product relation: $\mathcal{R}(\underline{X}_i, \underline{X}_j) = (\underline{X}_i - m)^\top (\underline{X}_j - m)$

- Linear mapping $\underline{X}' = \Phi(\underline{X}) = V^\top (\underline{X} - m)$.

- Euclidean scalar product matching:

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \left| (\underline{X}_i - m)^\top (\underline{X}_j - m) - \underline{X}'_i{}^\top \underline{X}'_j \right|^2$$

- Φ often defined only on $\mathcal{D} \dots$

MDS Heuristic

- Match the *scalar* products:

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \left| (\underline{X}_i - m)^\top (\underline{X}_j - m) - \underline{X}'_i{}^\top \underline{X}'_j \right|^2$$

- Linear method: $\underline{X}' = U^\top (\underline{X} - m)$ with U orthonormal

- **Beware:** \underline{X} can be unknown, only the scalar products are required!
- Resulting criterion: minimization in $U^\top (\underline{X}_i - m)$ of

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \left| (\underline{X}_i - m)^\top (\underline{X}_j - m) - (\underline{X}_i - m)^\top U U^\top (\underline{X}_j - m) \right|^2$$

without using explicitly \underline{X} in the algorithm...

- Explicit solution obtained through the eigendecomposition of the know Gram matrix $(\underline{X}_i - m)^\top (\underline{X}_j - m)$ by keeping only the d' largest eigenvalues.

- In this case, MDS yields the same result as the PCA (but with different inputs, distance between observation vs correlations)!
- **Explanation:** Same SVD problem up to a transposition:

- MDS

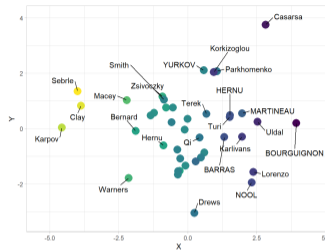
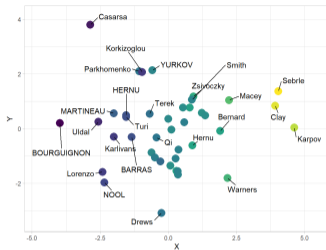
$$\underline{\bar{X}}_{(n)}^\top \underline{\bar{X}}_{(n)} \sim \underline{\bar{X}}_{(n)}^\top U U^\top \underline{\bar{X}}_{(n)}$$

- PCA

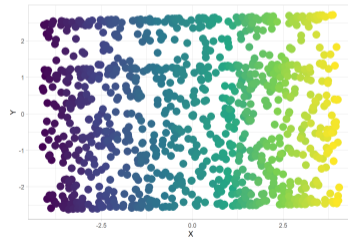
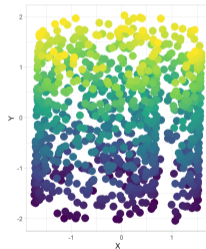
$$\underline{\bar{X}}_{(n)} \underline{\bar{X}}_{(n)}^\top \sim U^\top \underline{\bar{X}}_{(n)} \underline{\bar{X}}_{(n)}^\top U$$

- Complexity: PCA $O((n + d')d^2)$ vs MDS $O((d + d')n^2)$...

Decathlon



Swiss Roll



PCA

MDS

- Preserving the scalar products amounts to preserve the Euclidean distance.
- Easier **generalization** if we work in terms of distance!

Generalized MDS

- Generalized MDS:
 - Distance relation: $\mathcal{R}(\underline{X}_i, \underline{X}_j) = d(\underline{X}_i, \underline{X}_j)$
 - Linear mapping $\underline{X}' = \Phi(\underline{X}) = V^T(\underline{X} - m)$.
 - Euclidean matching:

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n |d(\underline{X}_i, \underline{X}_j) - d'(\underline{X}'_i, \underline{X}'_j)|^2$$

- Strong connection (but no equivalence) with MDS when $d(x, y) = \|x - y\|^2$!
- **Minimization:** Simple gradient descent can be used (can be stuck in local minima).

- MDS: equivalent to PCA (but more expensive) if $d(x, y) = \|x - y\|^2$!
- ISOMAP: use a *localized* distance instead to limit the influence of very far point.

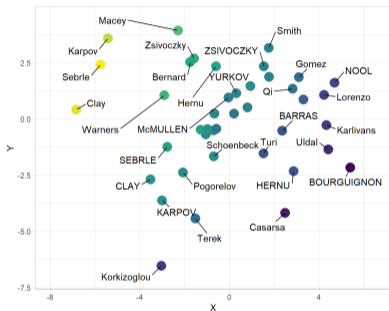
ISOMAP

- For each point \underline{X}_i , define a neighborhood \mathcal{N}_i (either by a distance or a number of points) and let

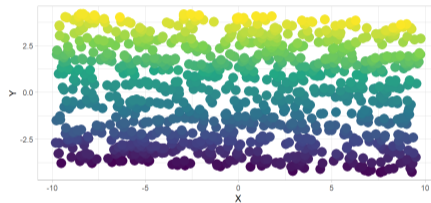
$$d_0(\underline{X}_i, \underline{X}_j) = \begin{cases} +\infty & \text{if } \underline{X}_j \notin \mathcal{N}_i \\ \|\underline{X}_i - \underline{X}_j\|^2 & \text{otherwise} \end{cases}$$

- Compute the shortest path distance for each pair.
- Use the MDS algorithm with this distance





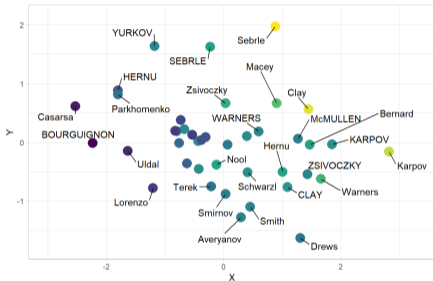
Decathlon



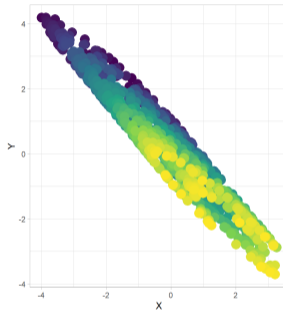
Swiss Roll

Random Projection Heuristic

- Draw at random d' unit vector (direction) U_i .
- Use $\underline{X}' = U^\top (\underline{X} - m)$ with $m = \frac{1}{n} \sum_{i=1}^n \underline{X}_i$
- **Property:** If \underline{X} lives in a space of dimension d'' , then, as soon as, $d' \sim d'' \log(d'')$,
$$\|\underline{X}_i - \underline{X}_j\|^2 \sim \frac{d}{d'} \|\underline{X}'_i - \underline{X}'_j\|^2$$
- Do not really use the data!



Decathlon



Swiss Roll

SNE heuristic

- From $\underline{X}_i \in \mathcal{X}$, construct a set of conditional probability:

$$P_{j|i} = \frac{e^{-\|\underline{X}_i - \underline{X}_j\|^2 / 2\sigma_i^2}}{\sum_{k \neq i} e^{-\|\underline{X}_i - \underline{X}_k\|^2 / 2\sigma_i^2}}$$

$$i \rightarrow P_i$$

$$P_{i|i} = 0$$

- Find \underline{X}'_i in $\mathbb{R}^{d'}$ such that the set of conditional probability:

$$Q_{j|i} = \frac{e^{-\|\underline{X}'_i - \underline{X}'_j\|^2 / 2\sigma_i^2}}{\sum_{k \neq i} e^{-\|\underline{X}'_i - \underline{X}'_k\|^2 / 2\sigma_i^2}}$$

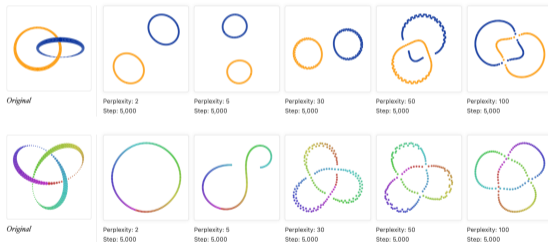
$$\sum_i d(P_i, Q_i)$$

$$Q_{i|i} = 0$$

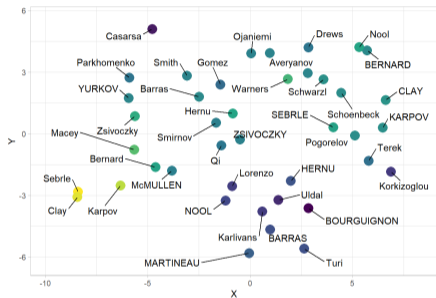
is close from P .

- t-SNE:** use a Student-t term $(1 + \|\underline{X}'_i - \underline{X}'_j\|^2)^{-1}$ for \underline{X}'_j
- Minimize the Kullback-Leibler divergence $(\sum_{i,j} P_{j|i} \log \frac{P_{j|i}}{Q_{j|i}})$ by a simple gradient descent (can be stuck in local minima).
- Parameters σ_i such that $H(P_i) = -\sum_{j=1}^n P_{j|i} \log P_{j|i} = \text{cst.}$

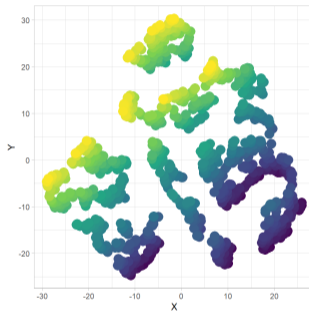
- Very successful/ powerful technique in practice
- Convergence may be long, unstable, or strongly depending on parameters.
- See this [distill post](#) for many impressive examples



Representation depending on t-SNE parameters



Decathlon



Swiss Roll

- Topological Data Analysis inspired.

Uniform Manifold Approximation and Projection

- Define a notion of asymmetric scaled local proximity between neighbors:
 - Compute the k -neighborhood of \underline{X}_i , its diameter σ_i and the distance ρ_i between \underline{X}_i and its nearest neighbor.
 - Define

$$w_i(\underline{X}_i, \underline{X}_j) = \begin{cases} e^{-(d(\underline{X}_i, \underline{X}_j) - \rho_i) / \sigma_i} & \text{for } \underline{X}_j \text{ in the } k\text{-neighborhood} \\ 0 & \text{otherwise} \end{cases}$$

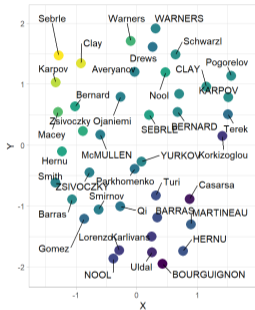
- Symmetrize into a *fuzzy* nearest neighbor criterion

$$w(\underline{X}_i, \underline{X}_j) = w_i(\underline{X}_i, \underline{X}_j) + w_j(\underline{X}_j, \underline{X}_i) - w_i(\underline{X}_i, \underline{X}_j)w_j(\underline{X}_j, \underline{X}_i)$$

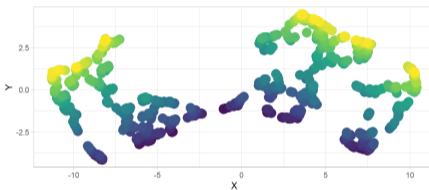
- Determine the points \underline{X}'_i in a low dimensional space such that

$$\sum_{i \neq j} w(\underline{X}_i, \underline{X}_j) \log \left(\frac{w(\underline{X}_i, \underline{X}_j)}{w'(\underline{X}'_i, \underline{X}'_j)} \right) + (1 - w(\underline{X}_i, \underline{X}_j)) \log \left(\frac{(1 - w(\underline{X}_i, \underline{X}_j))}{(1 - w'(\underline{X}'_i, \underline{X}'_j))} \right)$$

- Can be performed by local gradient descent.



Decathlon



Swiss Roll

Graph heuristic

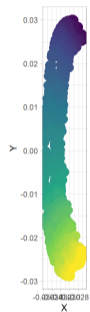
- Construct a graph with weighted edges $w_{i,j}$ measuring the *proximity* of \underline{X}_i and \underline{X}_j ($w_{i,j}$ large if close and 0 if there is no information).
- Find the points $\underline{X}'_j \in \mathbb{R}^{d'}$ minimizing

$$\frac{1}{n} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n w_{i,j} \|\underline{X}'_i - \underline{X}'_j\|^2$$

- Need of a constraint on the size of \underline{X}'_j ...
- Explicit solution through linear algebra: d' eigenvectors with smallest eigenvalues of the Laplacian of the graph $D - W$, where D is a diagonal matrix with $D_{i,i} = \sum_j w_{i,j}$.
- Variation on the definition of the Laplacian...



Decathlon



Swiss Roll

- 1 Unsupervised Learning?
- 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 3 **Dimension Reduction**
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - **Comparing Methods?**
- 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References

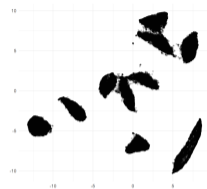
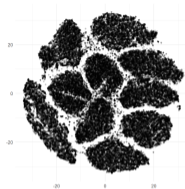
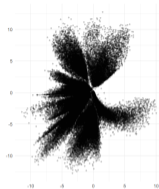
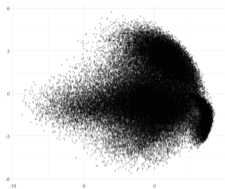
How to Compare Different Dimensionality Reduction Methods ?

- **Difficult!** Once again, the metric is very subjective.

However, a few possible attempts

- Did we preserve a lot of inertia with only a few directions?
- Do those directions *make sense* from an expert point of view?
- Do the low dimension representation *preserve* some important information?
- Are we better on **subsequent task**?

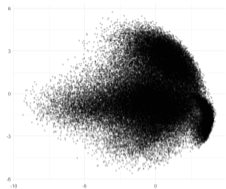
A Challenging Example: MNIST



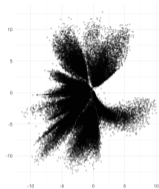
MNIST Dataset

- Images of 28×28 pixels.
- No label used!
- 4 different embeddings.

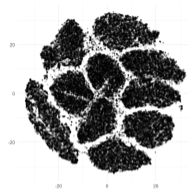
A Challenging Example: MNIST



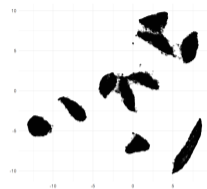
PCA



autoencoder



t-SNE

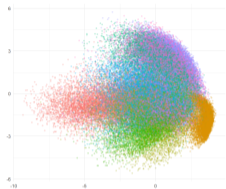


UMAP

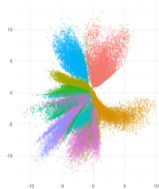
MNIST Dataset

- Images of 28×28 pixels.
- No label used!
- 4 different embeddings.

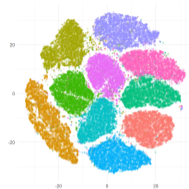
A Challenging Example: MNIST



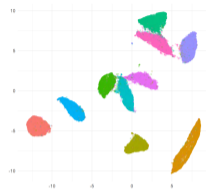
PCA



autoencoder



t-SNE

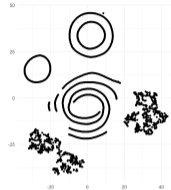
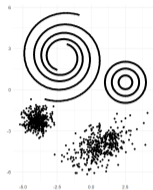


UMAP

MNIST Dataset

- Images of 28×28 pixels.
- No label used!
- 4 different embeddings.
- Quality evaluated by visualizing the true labels **not used to obtain the embeddings.**
- Only a few labels could have been used.

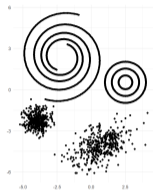
A Simpler Example: A 2D Set



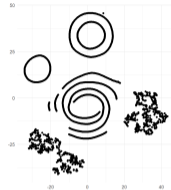
Cluster Dataset

- Set of points in 2D.
- No label used!
- 3 different embeddings.

A Simpler Example: A 2D Set



PCA



t-SNE

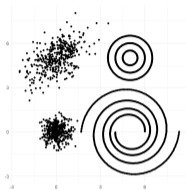


UMAP

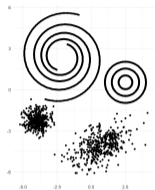
Cluster Dataset

- Set of points in 2D.
- No label used!
- 3 different embeddings.

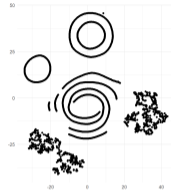
A Simpler Example: A 2D Set



Original



PCA



t-SNE



UMAP

Cluster Dataset

- Set of points in 2D.
- No label used!
- 3 different embeddings.

- Quality evaluated by stability...

- 1 Unsupervised Learning?
 - 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - 6 References

- **Training data** : $\mathcal{D} = \{\underline{X}_1, \dots, \underline{X}_n\} \in \mathcal{X}^n$ (i.i.d. $\sim \mathbb{P}$)
- Latent groups?

Clustering

- Construct a map f from \mathcal{D} to $\{1, \dots, K\}$ where K is a number of classes to be fixed:

$$f : \underline{X}_j \mapsto k_j$$

Motivations

- Interpretation of the groups
- Use of the groups in further processing
- Several strategies possible!
- Can use dimension reduction as a preprocessing.

- 1 Unsupervised Learning?
- 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 4 **Clustering**
 - **Prototype Approaches**
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References

Partition Heuristic

- Clustering is defined by a partition in K classes. . .
- that minimizes a homogeneity criterion.

K- Means

- Cluster k defined by a *center* μ_k .
- Each sample is associated to the closest center.
- Centers defined as the minimizer of $\sum_{i=1}^n \min_k \|\underline{X}_i - \mu_k\|^2$
- Iterative scheme (Lloyd):
 - Start by a (pseudo) random choice for the centers μ_k
 - Assign each samples to its nearby center
 - Replace the center of a cluster by the mean of its assigned samples.
 - Repeat the last two steps until convergence.

- Other schemes:
 - McQueen: modify the mean each time a sample is assigned to a new cluster.
 - Hartigan: modify the mean by removing the considered sample, assign it to the nearby center and recompute the new mean after assignment.

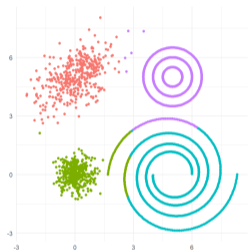
A good initialization is crucial!

- Initialize by samples.
- k-Mean++: try to take them as separated as possible.
- No guarantee to converge to a global optimum: repeat and keep the best result!
- Complexity : $O(n \times K \times T)$ where T is the number of steps in the algorithm.

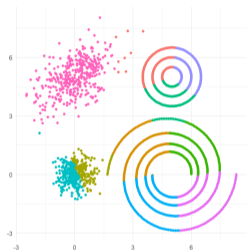
- k-Medoid: use a sample as a center
 - PAM: for a given cluster, use the sample that minimizes the intra distance (sum of the squared distance to the other points)
 - Approximate medoid: for a given cluster, assign the point that is the closest to the mean.

Complexity

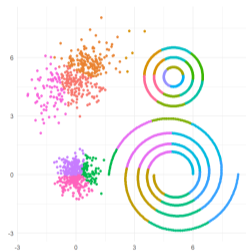
- PAM: $O(n^2 \times T)$ in the worst case!
- Approximate medoid: $O(n \times K \times T)$ where T is the number of steps in the algorithm.
- **Remark:** Any distance can be used... but the complexity of computing the centers can be very different.



$k = 4$



$k = 10$



$k = 10$

Model Heuristic

- Use a generative model of the data:

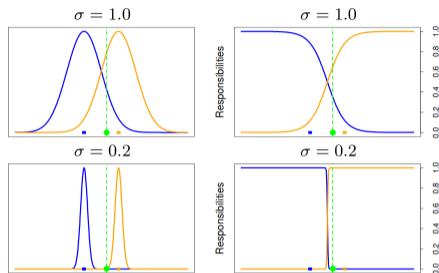
$$\mathbb{P}(\underline{X}) = \sum_{k=1}^K \pi_k \mathbb{P}_{\theta_k}(\underline{X}|k)$$

where π_k are proportions and $\mathbb{P}_{\theta}(\underline{X}|k)$ are parametric probability models.

- Estimate those parameters (often by a ML principle).
- Assign each observation to the class maximizing the a posteriori probability (obtained by Bayes formula)

$$\frac{\widehat{\pi}_k \mathbb{P}_{\widehat{\theta}_k}(\underline{X}|k)}{\sum_{k'=1}^K \widehat{\pi}_{k'} \mathbb{P}_{\widehat{\theta}_{k'}}(\underline{X}|k')}$$

- Link with Generative model in supervised classification!



A two class example

- A mixture $\pi_1 f_1(\underline{X}) + \pi_2 f_2(\underline{X})$
- and the posterior probability $\pi_i f_i(\underline{X}) / (\pi_1 f_1(\underline{X}) + \pi_2 f_2(\underline{X}))$
- Natural class assignment!

Sub-population estimation

- A mixture $\pi_1 f_1(\underline{X}) + \pi_2 f_2(\underline{X})$
- Two populations with a parametric distribution f_i .
- Most classical choice: Gaussian distribution

Gaussian Setting

- $\underline{X}_1, \dots, \underline{X}_n$ independent
- $\underline{X}_i \sim N(\mu_1, \sigma_1^2)$ with probability π_1 or $\underline{X}_i \sim N(\mu_2, \sigma_2^2)$ with probability π_2
- We don't know the parameters μ_i, σ_i, π_i .
- We don't know from which distribution each \underline{X}_i has been drawn.

Maximum Likelihood

- Density: $\pi_1 \Phi(\underline{X}, \mu_1, \sigma_1^2) + \pi_2 \Phi(\underline{X}, \mu_2, \sigma_2^2)$
- log-likelihood: $\mathcal{L}(\theta) = \sum_{i=1}^n \log (\pi_1 \Phi(\underline{X}_i, \mu_1, \sigma_1^2) + \pi_2 \Phi(\underline{X}_i, \mu_2, \sigma_2^2))$
- No straightforward way to optimize the parameters!

What if algorithm

- Assume we know from which distribution each sample has been sampled: $Z_i = 1$ if from f_1 and $Z_i = 0$ otherwise.
- log-likelihood: $\sum_{i=1}^n Z_i \log \Phi(\underline{X}_i, \mu_1, \sigma_1^2) + (1 - Z_i) \log \Phi(\underline{X}_i, \mu_2, \sigma_2^2)$
- Easy optimization... but the Z_i are unknown!

What if algorithm

- Assume we know from which distribution each sample has been sampled: $Z_i = 1$ if from f_1 and $Z_i = 0$ otherwise.
- log-likelihood:
$$\sum_{i=1}^n Z_i \log \Phi(\underline{X}_i, \mu_1, \sigma_1^2) + (1 - Z_i) \log \Phi(\underline{X}_i, \mu_2, \sigma_2^2)$$
- Easy optimization. . . but the Z_i are unknown!

Bootstrapping Idea

- Replace Z_i by its expectation given the current estimate.
- $\mathbb{E}[Z_i] = \mathbb{P}(Z_i = 1|\theta)$ (A posteriori probability)
- and iterate. . .
- Can be proved to be good idea!

EM Algorithm

- (Random) initialization: $\mu_i^0, \sigma_i^0, \pi_i^0$.
- Repeat:
 - Expectation (Current a posteriori probability):

$$\mathbb{E}_t[Z_i] = \mathbb{P}(Z_i = 1|\theta^t) = \frac{\pi_1^t \Phi(\underline{X}_i, \mu_1^t, (\sigma_1^t)^2)}{\pi_1^t \Phi(\underline{X}_i, \mu_1^t, (\sigma_1^t)^2) + \pi_2^t \Phi(\underline{X}_i, \mu_2^t, (\sigma_2^t)^2)}$$

- Maximization of

$$\sum_{i=1}^n \mathbb{E}_t[Z_i] \log \Phi(\underline{X}_i, \mu_1, \sigma_1^2) + \mathbb{E}_t[1 - Z_i] \log \Phi(\underline{X}_i, \mu_2, \sigma_2^2)$$

to obtain $\mu_i^{t+1}, \sigma_i^{t+1}, \pi_i^{t+1}$.

- Large choice of parametric models.

Gaussian Mixture Model

- Use

$$\mathbb{P}_{\theta_k}(\vec{X}|k) \sim \mathcal{N}(\mu_k, \Sigma_k)$$

with $\mathcal{N}(\mu, \Sigma)$ the Gaussian law of mean μ and covariance matrix Σ .

- Efficient optimization algorithm available (EM)
- Often some constraint on the covariance matrices: identical, with a similar structure. . .
- Strong connection with K -means when the covariance matrices are assumed to be the same multiple of the identity.

Probabilistic latent semantic analysis (PLSA)

- Documents described by their word counts w
- Model:

$$\mathbb{P}(w) = \sum_{k=1}^K \pi_k \mathbb{P}_{\theta_k}(w|k)$$

with k the (hidden) topic, π_k a topic probability and $\mathbb{P}_{\theta_k}(w|k)$ a multinomial law for a given topic.

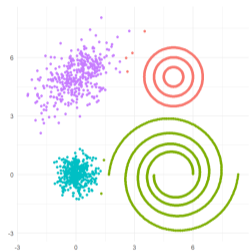
- Clustering according to

$$\mathbb{P}(k|w) = \frac{\hat{\pi}_k \mathbb{P}_{\hat{\theta}_k}(w|k)}{\sum_{k'} \hat{\pi}_{k'} \mathbb{P}_{\hat{\theta}_{k'}}(w|k')}$$

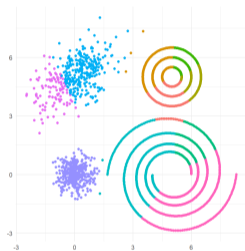
- Same idea than GMM!
- Bayesian variant called LDA.

Parametric Density Estimation Principle

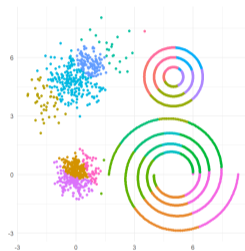
- Assign a probability of membership.
- Lots of theoretical studies. . .
- Model selection principle can be used to select K the number of classes (or rather to avoid using a nonsensical K . . .):
 - AIC / BIC / MDL penalization
 - Cross Validation is also possible!
- Complexity: $O(n \times K \times T)$



$k = 4$



$k = 10$



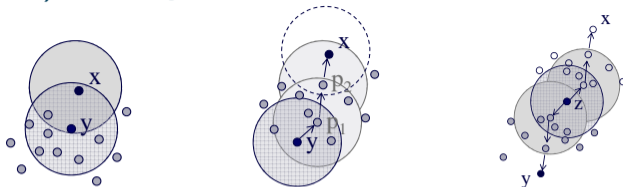
$k = 10$

- 1 Unsupervised Learning?
- 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 4 Clustering
 - Prototype Approaches
 - **Contiguity Approaches**
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References

Density Heuristic

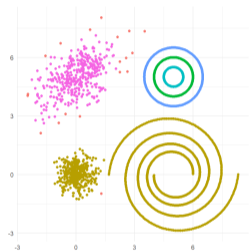
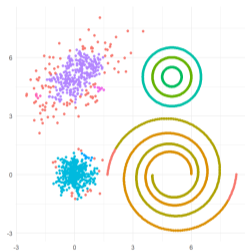
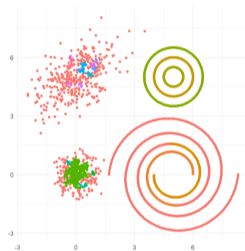
- Cluster are connected dense zone separated by low density zone.
- Not all points belong to a cluster.

- Basic bricks:
 - Estimate the density.
 - Find points with high densities.
 - Gather those points according to the density.
- Density estimation:
 - Classical kernel density estimators. . .
- Gathering:
 - Link points of high density and use the resulted component.
 - Move them toward top of density *hill* by following the gradient and gather all the points arriving at the same *summit*.



Examples

- DBSCAN: link point of high densities using a very simple kernel.
- PdfCLuster: find connected zone of high density.
- Mean-shift: move points toward top of density *hill* following an evolving kernel density estimate.
- Complexity: $O(n^2 \times T)$ in the worst case.
- Can be reduced to $O(n \log(n) T)$ if samples can be encoded in a tree structure (n-body problem type approximation).

 $\epsilon = .45$  $\epsilon = .2$  $\epsilon = .1$

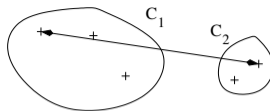
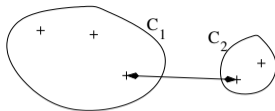
- 1 Unsupervised Learning?
- 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - **Agglomerative Approaches**
 - Other Approaches
 - Scalability
- 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References

Agglomerative Clustering Heuristic

- Start with very small clusters (a sample by cluster?)
 - Sequential merging of the most similar clusters. . .
 - according to some *greedy* criterion Δ .
-
- Generates a hierarchy of clustering instead of a single one.
 - Need to select the number of cluster afterwards.
 - Several choices for the merging criterion. . .
 - Examples:
 - Minimum Linkage: merge the closest cluster in term of the usual distance
 - Ward's criterion: merge the two clusters yielding the less inner inertia loss (k-means criterion)

Algorithm

- Start with $(\mathcal{C}_i^{(0)}) = (\{\underline{X}_i\})$ the collection of all singletons.
- At step s , we have $n - s$ clusters $(\mathcal{C}_i^{(s)})$:
 - Find the two most similar clusters according to a criterion Δ :
$$(i, i') = \underset{(j, j')}{\operatorname{argmin}} \Delta(\mathcal{C}_j^{(s)}, \mathcal{C}_{j'}^{(s)})$$
 - Merge $\mathcal{C}_i^{(s)}$ and $\mathcal{C}_{i'}^{(s)}$ into $\mathcal{C}_i^{(s+1)}$
 - Keep the $n - s - 2$ other clusters $\mathcal{C}_{i''}^{(s+1)} = \mathcal{C}_{i''}^{(s)}$
- Repeat until there is only one cluster.
- Complexity: $O(n^3)$ in general.
- Can be reduced to $O(n^2)$
 - if only a bounded number of merging is possible for a given cluster,
 - for the most classical distances by maintaining a nearest neighbors list.



Merging criterion based on the distance between points

- Minimum linkage:

$$\Delta(C_i, C_j) = \min_{\underline{X}_i \in C_i} \min_{\underline{X}_j \in C_j} d(\underline{X}_i, \underline{X}_j)$$

- Maximum linkage:

$$\Delta(C_i, C_j) = \max_{\underline{X}_i \in C_i} \max_{\underline{X}_j \in C_j} d(\underline{X}_i, \underline{X}_j)$$

- Average linkage:

$$\Delta(C_i, C_j) = \frac{1}{|C_i||C_j|} \sum_{\underline{X}_i \in C_i} \sum_{\underline{X}_j \in C_j} d(\underline{X}_i, \underline{X}_j)$$

- Clustering based on the proximity...

Merging criterion based on the inertia (distance to the mean)

- Ward's criterion:

$$\begin{aligned}\Delta(C_i, C_j) = & \sum_{\underline{X}_i \in C_i} \left(d^2(\underline{X}_i, \mu_{C_i \cup C_j}) - d^2(\underline{X}_i, \mu_{C_i}) \right) \\ & + \sum_{\underline{X}_j \in C_j} \left(d^2(\underline{X}_j, \mu_{C_i \cup C_j}) - d^2(\underline{X}_j, \mu_{C_j}) \right)\end{aligned}$$

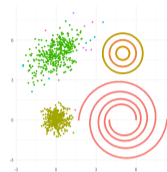
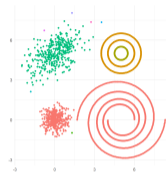
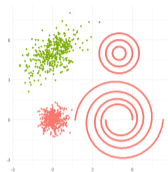
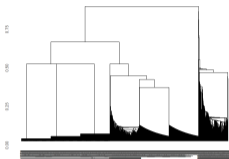
- If d is the Euclidean distance:

$$\Delta(C_i, C_j) = \frac{2|C_i||C_j|}{|C_i| + |C_j|} d^2(\mu_{C_i}, \mu_{C_j})$$

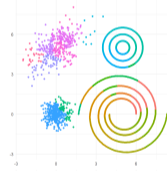
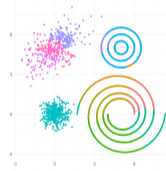
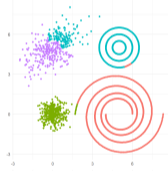
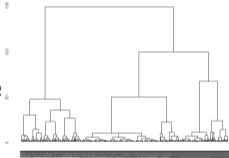
- Same criterion than in the k -means algorithm but greedy optimization.

Agglomerative Clustering

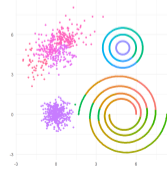
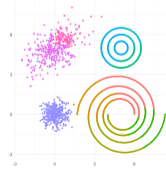
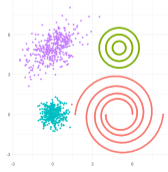
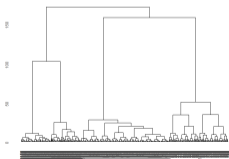
Single



Complete



Ward



Dendrogram

$k = 4$

$k = 10$

$k = 20$

- 1 Unsupervised Learning?
- 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - **Other Approaches**
 - Scalability
- 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References

Grid heuristic

- Split the space in pieces
- Group those of high density according to their proximity
- Similar to density based estimate (with partition based initial clustering)
- Space splitting can be fixed or adaptive to the data.
- Examples:
 - STING (Statistical Information Grid): Hierarchical tree construction plus DBSCAN type algorithm
 - AMR (Adaptive Mesh Refinement): Adaptive tree refinement plus k -means type assignment from high density leaves.
 - CLIQUE: Tensorial grid and 1D detection.
- Linked to Divisive clustering (DIANA)

Graph based

- Spectral clustering: dimension reduction + k-means.
- Message passing: iterative local algorithm.
- Graph cut: min/max flow.

- Kohonen Map,
- ...

- 1 Unsupervised Learning?
- 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - **Scalability**
- 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References

Large dataset issue

- When n is large, a $O(n^\alpha \log n)$ with $\alpha > 1$ is not acceptable!
- How to deal with such a situation?
- **Beware:** Computing all the pairwise distance requires $O(n^2)$ operations!

Ideas

- Sampling
- Online processing
- Simplification
- Parallelization

Sampling heuristic

- Use only a subsample to construct the clustering.
 - Assign the other points to the constructed clusters afterwards.
-
- Requires a clustering method that can assign new points (partition, model...)
 - Often repetition and choice of the best clustering
 - Example:
 - CLARA: K-medoid with sampling and repetition
 - Two-steps algorithm:
 - Generate a large number n' of clusters using a fast algorithm (with $n' \ll n$)
 - Cluster the clusters with a more accurate algorithm.

Online heuristic

- Modify the current clusters according to the value of a single observation.
- Requires compactly described clusters.
- Examples:
 - Add to an existing cluster (and modify it) if it is close enough and create a new cluster otherwise (k -means without reassignment)
 - Stochastic descent gradient (GMM)
- May leads to far from optimal clustering.

Simplification heuristic

- Simplify the algorithm to be more efficient at the cost of some precision.
- Algorithm dependent!
- Examples:
 - Replace groups of observation (preliminary cluster) by the (approximate) statistics.
 - Approximate the distances by cheaper ones.
 - Use n-body type techniques.

Parallelization heuristic

- Split the computation on several computers.
- Algorithm dependent!
- Examples:
 - Distance computation in k -means, parameter gradient in model based clustering
 - Grid density estimation, Space splitting strategies
- Classical batch sampling not easy to perform as partitions are not easily merged. . .

- 1 Unsupervised Learning?
- 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 5 **Generative Modeling**
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References

- **Training data** : $\mathcal{D} = \{(\underline{X}_1, \underline{Y}_1), \dots, (\underline{X}_n, \underline{Y}_n)\} \in (\mathcal{X} \times \mathcal{Y})^n$ (i.i.d. $\sim \mathbb{P}$)
- Same kind of data than for supervised learning if $\mathcal{Y} \neq \emptyset$.

Generative Modeling

- Construct a map G from the product of \mathcal{Y} and a randomness source Ω to \mathcal{X}

$$G : \mathcal{Y} \times \Omega \rightarrow \mathcal{X}$$

$$(Y, \omega) \mapsto X$$

- Unconditional model if $\mathcal{Y} = \emptyset$...

Motivation

- Generate plausible novel conditional samples based on a given dataset.

Sample Quality

- Related to the proximity between the law of $G(Y, \omega)$ and the law of $X|Y$.
- Most classical choice is the Kullback-Leibler divergence.

Ingredients

- Generator $G_\theta(Y, \omega)$ and cond. density prob. $P_\theta(X|Y)$ (Explicit vs implicit link)
- Simple / Complex / Approximate estimation...

Some Possible Choices

	Probabilistic model	Generator	Estimation
Base	Simple (parametric)	Explicit	Simple (ML)
Flow	Image of simple model	Explicit	Simple (ML)
Factorization	Factorization of simple model	Explicit	Simple (ML)
VAE	Simple model with latent var.	Explicit	Approximate (ML)
EBM	Arbitrary	Implicit (MCMC)	Complex (ML/score/discrim.)
Diffusion	Continuous noise	Implicit (MCMC)	Complex (score)
	Discrete Noise with latent var.	Explicit	Approximate (ML)
GAN	Implicit	Explicit	Complex (Discrimination)

- SOTA: Diffusion based approach!

$$\tilde{X} = G(Y, \omega) \quad ?$$

- Small abuse of notations...
- More an algorithm than a map!

Generators

- One step: $\omega \sim \tilde{Q}(\cdot|Y)$ and $\tilde{X} = G(Y, \omega)$.
- Several steps:
 - $\omega_0 \sim \tilde{Q}_0(\cdot|Y)$ and $\tilde{X}_0 = G_0(\omega_0)$
 - $\omega_{t+1} \sim \tilde{Q}_{t+1}(\cdot|Y, \tilde{X}_t)$ and $\tilde{X}_{t+1} = G_{t+1}(Y, \tilde{X}_k, \omega_{t+1})$
- Fixed or variable number of steps.
- Fixed or variable dimension for \tilde{X}_t and ω_t ...
- \tilde{Q} (or \tilde{Q}_t) should be easy to sample.
- Most of the time, parametric representations for \tilde{Q} (or \tilde{Q}_t) and G (or G_t).

- 1 Unsupervised Learning?
 - 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - 5 **Generative Modeling**
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - 6 References

$$X \sim P \text{ with } dP(x) = p(x)d\lambda \longrightarrow \tilde{X} \sim \tilde{P} \text{ with } d\tilde{P}(x) = \tilde{p}(x)d\lambda$$

Heuristic

- Estimate p by \tilde{p} from an i.i.d. sample X_1, \dots, X_n .
- Simulate \tilde{X} having a law \tilde{P} .
- By construction, if \tilde{p} is *close* from p , the law of \tilde{X} will be close from the law of X .

Issue: How to do it?

- How to estimate \tilde{p} ? Parametric, non-parametric? Maximum likelihood? Other criteria?
- How to simulate \tilde{P} ? Parametric? One-step? Multi-step? Iterative?

$X \sim P(\cdot)$ with $dP(x) = p(x)d\lambda \longrightarrow \tilde{X} \sim \tilde{P}_\theta$ with $d\tilde{P}_\theta(x) = \tilde{p}_\theta(x)d\lambda$

Maximum Likelihood Approach

- Select a family \tilde{P} and estimate p by \tilde{p}_θ from an i.i.d. sample X_1, \dots, X_n .
- Simulate \tilde{X} having a law \tilde{P}_θ .
- By construction, if \tilde{p}_θ is *close* from p , the law of \tilde{X} will be close from the law of X .

Issue: How to do it?

- Which family \tilde{P} ?
- How to simulate \tilde{P}_θ ? Parametric? Iterative?
- Corresponds to $\omega \sim \tilde{P}_\theta$ and $\tilde{X} = G(\omega) = \omega$

$$X|Y \sim P(\cdot | Y) \text{ with } dP(x|Y) = p(x|Y)d\lambda$$

$$\longrightarrow \tilde{X}|Y \sim \tilde{P}(\cdot | Y) \text{ with } d\tilde{P}(x|Y) = \tilde{p}(x|Y)d\lambda$$

Heuristic

- Estimate p by \tilde{p} from an i.i.d. sample $(X_1, Y_1), \dots, (X_n, Y_n)$.
- Simulate $\tilde{X}|Y$ having a law $\tilde{P}(\cdot | Y)$.
- By construction, if \tilde{p} is *close* from p , the law of $\tilde{X}|Y$ will be close from the law of $X|Y$.

Issue: How to do it?

- How to estimate \tilde{p} ? Parametric, non-parametric? Maximum likelihood? Other criteria?
- How to simulate \tilde{P} ? Parametric? One-step? Multi-step? Iterative?

$$X|Y \sim P(\cdot|Y) \text{ with } dP(x|Y) = p(x|Y)d\lambda$$
$$\longrightarrow \tilde{X}|Y \sim \tilde{P}_{\tilde{\theta}(Y)} \text{ with } d\tilde{P}_{\tilde{\theta}(Y)}(x) = \tilde{p}_{\tilde{\theta}(Y)}(x)d\lambda$$

Maximum Likelihood Approach

- Select a family \tilde{P} and estimate p by $\tilde{p}_{\tilde{\theta}}$ from an i.i.d. sample $(X_1, Y_1), \dots, (X_n, Y_n)$ where $\tilde{\theta}$ is now a function of Y .
- Simulate $\tilde{X}|Y$ having a law $\tilde{P}_{\tilde{\theta}(Y)}$
- If $\tilde{p}_{\tilde{\theta}}$ is close from p , the law of $\tilde{X}|Y$ will be close from the law of $X|Y$.

Issue: How to do it?

- Which family \tilde{P} ? Which function family for $\tilde{\theta}$?
- How to simulate $\tilde{P}_{\tilde{\theta}(Y)}$? Parametric? Iterative?
- Corresponds to $\omega \sim \tilde{Q}(\cdot|Y) = \tilde{P}_{\tilde{\theta}(Y)}$ and $\tilde{X} = G(Y, \omega) = \omega$

$$\omega \sim \tilde{Q}_{\tilde{\theta}(Y)} \sim \tilde{q}_{\tilde{\theta}(Y)}(x) d\lambda \quad \text{and} \quad \tilde{X}|Y = G(Y, \omega) = \omega$$

Estimation

- By construction,

$$dP(\tilde{X}|Y) = \tilde{q}_{\tilde{\theta}(Y)}(x) d\lambda$$

- Maximum Likelihood approach:

$$\tilde{\theta} = \operatorname{argmax}_{\theta} \sum_{i=1}^n \log \tilde{q}_{\tilde{\theta}(Y_i)}(X_i)$$

Simulation

- \tilde{P} has been chosen so that this distribution is easy to sample...
- Possible families: Gaussian, Multinomial, Exponential model...
- Possible parametrizations for $\tilde{\theta}$: linear, neural network...
- Limited expressivity!

$\omega \sim \tilde{Q}_{\tilde{\theta}(Y)} \sim \tilde{q}_{\tilde{\theta}(Y)}(x)d\lambda$ and $\tilde{X}|Y = G(\omega)$ with a given G invertible.

Estimation

- By construction,

$$d\tilde{P}(G^{-1}(\tilde{X})|Y) = \tilde{q}_{\tilde{\theta}(Y)}(G^{-1}(x))d\lambda$$

- Maximum Likelihood approach:

$$\tilde{\theta} = \operatorname{argmax}_{\theta} \sum_{i=1}^n \log \tilde{q}_{\tilde{\theta}(Y_i)}(G^{-1}(X_i)) + \underbrace{\dots}_{??}$$

Simulation

- \tilde{Q} has been chosen so that this distribution is easy to sample...
- Possible transform G : Change of basis, known transform...

$\omega \sim \tilde{Q}_{\tilde{\theta}(Y)} = \tilde{q}_{\tilde{\theta}(Y)}(x) d\lambda$ and $\tilde{X}|Y = G_{\tilde{\theta}_G(Y)}(\omega)$ with G_{θ} invertible.

Estimation

- By construction,

$$d\tilde{P}(\tilde{X}|Y) = |\text{Jac} G_{\tilde{\theta}_G(Y)}^{-1}(x)| \tilde{q}_{\tilde{\theta}(Y)}(G_{\tilde{\theta}_G(Y)}^{-1}(x)) d\lambda$$

where $\text{Jac} G_{\tilde{\theta}_G(Y)}^{-1}(x)$ is the Jacobian of $G_{\tilde{\theta}_G(Y)}^{-1}$ at x

- Maximum Likelihood approach:

$$\tilde{\theta}, \tilde{\theta}_G = \underset{\theta, \theta_G}{\text{argmax}} \sum_{i=1}^n \left(\log |\text{Jac} G_{\theta_G(Y)}^{-1}(x)| + \log \tilde{q}_{\theta(Y)}(G_{\theta_G(Y)}^{-1}(X_i)) \right)$$

Simulation

- \tilde{Q} has been chosen so that this distribution is easy to sample...
- Often, in practice, $\tilde{\theta}(Y)$ is independent of Y ...
- Main issue: G_{θ} , its inverse and its Jacobian should be easy to compute.

$G_\theta?$

- Main issue: G_θ , its inverse and its Jacobian should be easy to compute.

Flow Models

- Composition

$$G_\theta = G_{\theta_T} \circ G_{\theta_{T-1}} \circ G_{\theta_1} \circ G_{\theta_0}$$

$$|\text{Jac}G_\theta^{-1}| = \prod |\text{Jac}G_{\theta_i}^{-1}|$$

- Real NVP

$$G_\theta(x) = \begin{pmatrix} x_1 \\ \vdots \\ x_{d'} \\ x_{d'+1}e^{s_{d'+1}(x_1, \dots, d')} + t_d(x_1, \dots, d') \\ \vdots \\ x_d e^{s_d(x_1, \dots, d')} + t_d(x_1, \dots, d') \end{pmatrix} \rightarrow G_\theta^{-1}(x) = \begin{pmatrix} x_1 \\ \vdots \\ x_{d'} \\ (x_{d'+1} - t_d(x_1, \dots, d'))e^{-s_{d'+1}(x_1, \dots, d')} \\ \vdots \\ (x_d - t_d(x_1, \dots, d'))e^{-s_d(x_1, \dots, d')} \end{pmatrix} \rightarrow |\text{Jac}G(x)^{-1}| = \prod_{d''=d'+1}^d e^{-s_{d''}(x_1, \dots, d')}$$

- Combined with permutation along dimension or invertible transform across dimension.
- Not that much flexibility...

$$\omega_0 \sim \tilde{Q}_0(\cdot|Y) \text{ and } \tilde{X}_0 = G_0(\omega_0)$$

$$\omega_{t+1} \sim \tilde{Q}_{t+1}(\cdot|Y, (\tilde{X}_l)_{l \leq t}) \text{ and } \tilde{X}_{t+1} = G_{t+1}(Y, (\tilde{X}_l)_{l \leq t}, \omega_{t+1})$$

$$\tilde{X} = (\tilde{X}_0, \dots, \tilde{X}_{d-1})$$

Factorization

- Amounts to use a factorized representation

$$\tilde{P}(\tilde{X}|Y) = \prod_{0 \leq k < d} \tilde{P}(\tilde{X}_k|Y, (\tilde{X}_l)_{l < k})$$

- \tilde{Q}_t and G_k can be chosen as in the plain density estimation case as the $X_i^{(t)}$ are observed.

Estimation

- d generative models to estimate instead of one.

- Simple generator by construction.

$$\omega_{t+1} \sim \tilde{Q}(\cdot | Y, (\tilde{X}_l)_{t \geq l \geq t-o}) \text{ and } \tilde{X}_{t+1} = G(Y, (\tilde{X}_l)_{t \geq l \geq t-o}, \omega_{t+1})$$
$$\tilde{X} = (\tilde{X}_0, \dots, \tilde{X}_{d-1})$$

Sequence and Markov Models

- Sequence: sequence of *similar* objects with a translation invariant structure.
 - Translation invariant probability model of finite order (memory) o .
 - Requires an initial padding of the sequence.
-
- Faster training as the parameters are shared for all t .
 - Model used in Text Generation!



- 1 Unsupervised Learning?
- 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 5 **Generative Modeling**
 - (Plain) Parametric Density Estimation
 - **Latent Variables**
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 References

$$\omega_0 \sim \tilde{Q}_0(\cdot | Y) \text{ and } \tilde{X}_0 = G_0(\omega_0)$$

$$\omega_1 \sim \tilde{Q}_1(\cdot | Y, X_0) \text{ and } \tilde{X}_1 = G_1(Y, \omega_0)$$

$$\tilde{X} = \tilde{X}_1$$

- Most classical example:
 - Gaussian Mixture Model with $\tilde{X}_0 = \omega_0 \sim \mathcal{M}(\pi)$ and $\tilde{X} = \omega_1 \sim \mathcal{N}(\mu_{\tilde{X}_0}, \Sigma_{\tilde{X}_0})$.

Estimation

- Still a factorized representation

$$\tilde{P}(\tilde{X}_1, \tilde{X}_0 | Y) = \tilde{P}_0(\tilde{X}_0 | Y) \tilde{P}_1(\tilde{X}_1 | Y, \tilde{X}_0)$$

but only \tilde{X}_1 is observed.

- **Much more complex estimation!**
- Simple generator by construction provided that the \tilde{Q}_t are easy to simulate.

$$\begin{aligned} \log \tilde{p}(\tilde{X}|Y) &= \log \mathbb{E}_{\tilde{P}(\tilde{X}_0|Y, \tilde{X})} [\tilde{p}(\tilde{X}, \tilde{X}_0|Y)] \\ &= \sup_{R(\cdot|Y, \tilde{X})} \mathbb{E}_{R(\cdot|Y, \tilde{X})} [\log \tilde{p}(\tilde{X}, \tilde{X}_0|Y) - \log r(\tilde{X}_0|Y, \tilde{X})] \end{aligned}$$

- Need to integrate over \tilde{X}_0 using the conditional law $\tilde{P}(\tilde{X}_0|Y, \tilde{X})$, which may be hard to compute.

Evidence Lower Bound

- Using $\log \tilde{p}(\tilde{X}|Y) = \mathbb{E}_{R(\cdot|Y, \tilde{X})} [\log (\tilde{p}(\tilde{X}, \tilde{X}_0|Y) / \tilde{p}(\tilde{X}_0|Y, \tilde{X}))]$,

$$\log \tilde{p}(\tilde{X}|Y) = \mathbb{E}_{R(\cdot|Y, \tilde{X})} [\log \tilde{p}(\tilde{X}, \tilde{X}_0|Y) - \log r(\tilde{X}_0|Y, \tilde{X}) - \text{KL}_{\tilde{X}_0}(R(\tilde{X}_0|Y, \tilde{X}), \tilde{P}(\tilde{X}_0|Y, \tilde{X}))]$$
- Equality is obtained for $R(\cdot|Y, \tilde{X}) = \tilde{P}(\tilde{X}_0|Y, \tilde{X})$.

- Maximization over \tilde{P} and R instead of only over \tilde{P} ...

$$\begin{aligned}
 \sup_{\tilde{P}} \mathbb{E}_{\tilde{X}, Y} [\log \tilde{p}(\tilde{X} | Y)] &= \sup_{\tilde{P}, R} \mathbb{E}_{\tilde{X}, Y, \tilde{X}_0 \sim R(\cdot | Y, \tilde{X})} [\log \tilde{p}(\tilde{X}, \tilde{X}_0 | Y) - \log r(\tilde{X}_0 | Y, \tilde{X})] \\
 &= \sup_{\tilde{P}, R} \mathbb{E}_{\tilde{X}, Y, \tilde{X}_0 \sim R(\cdot | Y, \tilde{X})} [\log \tilde{p}(\tilde{X} | Y, \tilde{X}_0)] \\
 &\quad + \underbrace{\mathbb{E}_{\tilde{X}, Y, \tilde{X}_0 \sim R(\cdot | Y, \tilde{X})} [\log \tilde{p}(\tilde{X}_0 | Y) - \log r(\tilde{X}_0 | Y, \tilde{X})]}_{\mathbb{E}_{\tilde{X}, Y} [\text{KL}(R(\cdot | Y, \tilde{X}), \tilde{P}(\tilde{X}_0 | Y))]}
 \end{aligned}$$

- Parametric models for $\tilde{P}(X_0 | Y)$, $\tilde{P}(X | Y, X_0)$ and $R(\tilde{X}_0 | \tilde{X}, Y)$.

Stochastic Gradient Descent

- Sampling on $(X, Y, \tilde{X}_0 \sim R)$ for $\mathbb{E}_{\tilde{X}, Y, \tilde{X}_0 \sim R(\cdot | Y, \tilde{X})} [\nabla \log \tilde{p}(\tilde{X} | Y, \tilde{X}_0)]$
- Sampling on (X, Y) for $\mathbb{E}_{\tilde{X}, Y} [\nabla \text{KL}(R(\cdot | Y, \tilde{X}), \tilde{P}(\tilde{X}_0 | Y))]$ if closed formula.
- Reparametrization trick for the second term otherwise...

Reparametrization Trick

$$\mathbb{E}_X[\beta(x)]$$

$$X = G(\omega) \quad \text{Generative Modeling} \\ \mathbb{E}_Z(\beta(Z))$$



$$Z = G(\omega) \text{ with } \omega \sim Q(\cdot) \text{ fixed} \longrightarrow \nabla \mathbb{E}_Z[F(Z)] = \nabla \mathbb{E}_\omega[F(G(\omega))] = \mathbb{E}_\omega[\nabla(F \circ G)(\omega)]$$

Reparametrization Trick

- Define a random variable Z as the image by a parametric map G of a random variable ω of fixed distribution Q .
- Most classical case: Gaussian...
- Allow to compute the derivative the expectation of a function of Z through a sampling of ω .

- Application for ELBO:

- $\tilde{X}_0 = G_R(Y, \tilde{X}, \omega_R)$ with $\omega_R \sim Q(\cdot | Y, \tilde{X})$ a fixed probability law.
- Sampling on ω to approximate:

$$\begin{aligned} \nabla \mathbb{E}_{\tilde{X}, Y, \tilde{X}_0 \sim R(\cdot | Y, \tilde{X})} \left[\log \tilde{p}(\tilde{X}_0 | Y) - \log r(\tilde{X}_0 | Y, \tilde{X}) \right] \\ = \mathbb{E}_{\tilde{X}, Y, \omega_0 \sim Q(\cdot | Y, \tilde{X})} \left[\nabla \log \tilde{p}(G_R(Y, \tilde{X}, \omega_0) | Y) - \nabla \log r(G_R(Y, \tilde{X}, \omega_0) | Y, \tilde{X}) \right] \end{aligned}$$

$$Z \sim \mathcal{N}(\mu, \sigma^2)$$

$$Z = \mu + \sigma \tilde{Z}$$

$$\tilde{Z} \sim \mathcal{N}(0, 1)$$

$$\mathbb{E}_Z[\ell(Z)] = \mathbb{E}_{\tilde{Z}}[\ell(\mu + \sigma \tilde{Z})]$$

$$\mathbb{E}[g(w)] \approx \frac{1}{m} \sum g(w_i)$$

$$\nabla_{\mu} \mathbb{E}_Z[\ell(Z)] = \nabla_{\mu} \mathbb{E}_{\tilde{Z}}[\ell(\mu + \sigma \tilde{Z})]$$

$$\nabla_{\mu} \left[\int \ell(x) \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx \right]$$

$$\mathbb{E}_{\tilde{Z}} \left[\nabla_{\mu} \ell(\mu + \sigma \tilde{Z}) \right]$$

$$\nabla_{\mu} = \int + 2(x-\mu) \ell(x) \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$$

$$\approx \frac{1}{m} \sum \nabla_{\mu} \ell(\mu + \sigma \tilde{Z}_i)$$

Generation: $\tilde{X}_0 \sim \tilde{P}(\cdot|Y) \xrightarrow{\text{decoder}} \tilde{X} \sim \tilde{P}(\cdot|Y, X_0)$

Training: $X \sim P(\cdot|Y) \xrightarrow{\text{encoder}} X_0 \sim R(\cdot|Y, X) \xrightarrow{\text{decoder}} \tilde{X} \sim \tilde{P}(\cdot|Y, X_0)$

Variational Auto Encoder

- Training structure similar to classical autoencoder... but matching on distributions rather than samples.
- Encoder interpretation of the approximate posterior $R(\cdot|Y, X)$.
- Implicit *low* dimension for X_0 .

$$\omega_0 \sim \tilde{Q}_0(\cdot|Y) \text{ and } \tilde{X}_0 = G_0(\omega_0)$$

$$\omega_{t+1} \sim \tilde{Q}_{t+1}(\cdot|Y, X_k) \text{ and } \tilde{X}_{t+1} = G_{t+1}(Y, \tilde{X}_t, \omega_{t+1})$$

$$\tilde{X} = \tilde{X}_T$$

Latent Variables

- Deeper hierarchy is possible. . .
- ELBO scheme still applicable using *decoders* R_i

$$R_i(\tilde{X}_i|Y, \tilde{X}_{i+1}) \simeq \tilde{P}(\tilde{X}_i|Y, \tilde{X}_{i+1})$$

- 1 Unsupervised Learning?
 - 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - 5 **Generative Modeling**
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - **Approximate Simulation**
 - Diffusion Model
 - Generative Adversarial Network
 - 6 References

$$d\tilde{P}(\tilde{X}|Y) \propto e^{u(\tilde{X}, Y)} d\lambda$$

$$\longrightarrow \omega_{t+1} \sim \tilde{Q}_u(\cdot|Y, \tilde{X}_t) \text{ and } \tilde{X}_{t+1} = G_u(Y, X_t, \omega_{t+1})$$

$$\tilde{X} \simeq \lim \tilde{X}_t$$

- Explicit conditional density model up to normalizing constant

$$Z(Y) = \int e^{u(x, Y)} d\lambda(x)$$

Simulation

- Several MCMC schemes to simulate the law without knowing $Z(Y)$

Estimation

- Not so easy as $Z(Y)$ depends a lot on u itself.

$$X_{t+1/2} \sim \tilde{Q}_u(\cdot | Y, \tilde{X}_t)$$

$$X_{t+1} = \begin{cases} X_{t+1/2} & \text{with proba } \alpha_t \\ X_t & \text{with proba } 1 - \alpha_t \end{cases} \quad \text{with } \alpha_t = \min \left(1, \frac{e^{u(X_{t+1/2}, Y)} \tilde{Q}_u(\tilde{X}_t | Y, \tilde{X}_{t+1/2})}{e^{u(X_t, Y)} \tilde{Q}_u(\tilde{X}_{t+1/2} | Y, \tilde{X}_t)} \right)$$

Metropolis Hastings

- Most classical algorithm.
- Convergence guarantee under reversibility of the proposal.
- Main issue is the choice of this proposal \tilde{Q} .
- Many enhanced versions exist!

Monte Carlo Markov Chain

$$\tilde{X}_{t+1/2} = X_t + \gamma_t \nabla_{\tilde{X}} u(\tilde{X}_t, Y) + \sqrt{2\gamma_t} \omega_t \text{ with } \omega_t \sim \mathcal{N}(0, \text{Id})$$

$$\tilde{X}_{t+1} = \begin{cases} \tilde{X}_{t+1/2} & \text{with proba } \alpha_t \\ \tilde{X}_t & \text{with proba } 1 - \alpha_t \end{cases} \text{ with } \alpha_t = \min \left(1, \frac{e^{u(\tilde{X}_{t+1/2}, Y)} \tilde{Q}_u(\tilde{X}_t | Y, \tilde{X}_{t+1/2})}{e^{u(\tilde{X}_t, Y)} \tilde{Q}_u(\tilde{X}_{t+1/2} | Y, \tilde{X}_t)} \right)$$

- Connection with a SDE:

$$\frac{d\tilde{X}}{dt} = \nabla_{\tilde{X}} u(\tilde{X}, Y) + \sqrt{2} dB_t$$

where B_t is a Brownian Motion.

Langevin

- If $\gamma_t = \gamma$, Metropolis-Hasting algorithm.
- With $\tilde{X}_{t+1} = \tilde{X}_{t+1/2}$, convergence toward an approximation of the law.
- Connection with SGD with decaying α_t

$$X|Y \sim P(\cdot|Y) \longrightarrow \tilde{X}|Y \sim \tilde{P}(\cdot|Y) \text{ with } d\tilde{P}(x|Y) = \tilde{p}(x|Y)d\lambda \propto e^{u(x,Y)}d\lambda$$

- Intractable log-likelihood:

$$\log \tilde{p}(\tilde{x}|Y) = u(\tilde{x}, Y) - \log Z(u)$$

Estimation

- Contrastive: simulate some \tilde{P} at each step and use

$$\nabla \log \tilde{p}(\tilde{x}|Y) = \nabla u(\tilde{x}, Y) - \nabla \log Z(u) = \nabla_{\tilde{x}} u(\tilde{x}, Y) - \mathbb{E}_{\tilde{P}} [\nabla_{\tilde{x}} u(\tilde{X}, Y)]$$

- Noise contrastive: learn to discriminate $W = X$ from $W = X' \sim R(\cdot|Y) \sim e^{r(x,Y)}d\lambda$ with the parametric approximation

$$\mathbb{P}(W = X|Y) \simeq \frac{e^{u(x,Y)}}{e^{u(x,Y)} + Z(Y)e^{r(x,Y)}}$$

- Score based: learn directly $s(\cdot|Y) = \nabla_{\tilde{x}} u(\cdot|Y) \simeq \nabla_X \log p(\cdot|Y)$.

$$\mathbb{E} \left[\|\nabla_X \log p(X|Y) - s(X|Y)\|^2 \right] = \mathbb{E} \left[\frac{1}{2} \|s(x|Y)\|^2 + \text{tr} \nabla_X s(X, Y) \right] + \text{cst.}$$

Score Based Method

- Non trivial formula based on partial integration.
- Hard to use in high dimension

$$\begin{aligned} X_\sigma = X + \sigma\epsilon &\longrightarrow \mathbb{E} \left[\|\nabla_X \log p_\sigma(X_\sigma|Y) - s_\sigma(W, Y)\|^2 \right] \\ &= \mathbb{E} \left[\|\nabla_X \log p_\sigma(X_\sigma|X, Y) - s_\sigma(X_\sigma, Y)\|^2 \right] + \text{cst.} \end{aligned}$$

Noisy Score

- Connection to denoising through Tweedie formula

$$\mathbb{E}[X|X_\sigma] = X_\sigma + \sigma^2 \nabla_X \log p_\sigma(X_\sigma|X, Y) \text{ and thus } s_\sigma(X_\sigma, Y) \simeq \frac{\mathbb{E}[X|X_\sigma] - X_\sigma}{\sigma^2}$$

$$\tilde{X} \sim e^{u(X,Y)} d\lambda \longrightarrow \tilde{X}_T \sim e^{\frac{1}{T}u(X,Y)}$$



Annealing

- Simulate a sequence of \tilde{X}_T starting with T large and decaying to 1.

$$\begin{aligned} X_\sigma = X + \sigma\epsilon &\longrightarrow \mathbb{E} \left[\|\nabla_X \log p_\sigma(X_\sigma | Y) - s_\sigma(W, Y)\|^2 \right] \\ &= \mathbb{E} \left[\|\nabla_X \log p_\sigma(X_\sigma | X, Y) - s_\sigma(X_\sigma, Y)\|^2 \right] + \text{cst.} \end{aligned}$$

Noisy Score

- Simulate a noisy sequence of \tilde{X}_σ with σ decaying to 0.

- 1 Unsupervised Learning?
 - 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - 5 **Generative Modeling**
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - **Diffusion Model**
 - Generative Adversarial Network
 - 6 References

Forward: $X_{t+1} = X_t + \sigma_t \omega_t \rightarrow X_{t+1} | X_0 \sim N(X_0, s_t^2 = \sum_{t' \leq t} \sigma_{t'}^2)$

Reverse: $\tilde{X}_T \sim N(0, s_T^2) \rightarrow \tilde{X}_{t-1} = \tilde{X}_t + \gamma_t u_{s_t^2}(X_t) + \sqrt{2\gamma_t} \omega_t$

Noisy Model

- Construct a sequence of noisy version X_t .
- Each X_t is easily sampled from X_0 so that the scores $u_{s_t^2}$ can be estimated.
- Use a simplistic sequential Langevin approach to obtain $\tilde{X}_0 \sim \tilde{P}(X)$ from $\tilde{X}_T \sim N(0, s_T^2)$.
- Lot of approximations everywhere.
- Dependency on Y has been removed for sake of simplicity.

Forward:
$$X_{t+\delta_t} = (1 + \alpha_t \delta_t) X_t + \sqrt{2\beta_t \delta_t} \omega_t$$
$$\longrightarrow dX(t) = \alpha(t) X(t) dt + \sqrt{2\beta(t)} dB(t)$$

Forward diffusion

- Generalization of noisy model:

$$X(t)|X(0) = N \left(X(0) \exp \int_0^t \alpha(u) du, \int_0^t 2\beta(u) \exp \left(\int_u^t \alpha(v) dv \right) du \right)$$

Reverse:
$$dX(t) = (-2\beta(t) \nabla_X \log P(X, t) - \alpha(t) X(t)) \bar{dt} + \sqrt{2\beta(t)} \bar{dB}(t)$$
$$\longrightarrow X_{t-\delta_t} = (1 - \alpha_t \delta_t) X_t + 2\beta_t \nabla_X \log p(X, t) \delta_t + \sqrt{2\beta_t \delta_t} \omega_t$$

Reverse diffusion

- Allow to sample back in time $X_t | X_T$.
- Quite involved derivation... but Langevin type scheme in the end.

$$\alpha_t = 0 \rightarrow \alpha(t) = 0 \rightarrow X(t)|X(0) = N\left(X(0), 2 \int_0^t \beta(u) du\right)$$

Noise Conditioned Score

- Direct extension of noisy model.
- Better numerical scheme but numerical explosion for $X(t)$.

$$(1 + \alpha_t \delta_t) = \sqrt{1 - 2\beta_t \delta_t} \simeq 1 - \beta_t \delta_t$$
$$\rightarrow X(t)|X(0) = N\left(X(0) \int_0^t e^{-\beta(u)} du, 2 \left(1 - \int_0^t e^{-\beta(u)}\right)\right)$$

Denoising Diffusion Probabilistic Model

- Explicit decay of the dependency on $P(X)$ and control on the variance.
- Better numerical results.
- Scores $\nabla_X \log p(X, t)$ estimated using the denoising trick as $X(t)|X(0)$ is explicit.
- Choice of $\beta(t)$ has a numerical impact.

Forward (SDE): $dX(t) = \alpha(t)X(t)dt + \sqrt{2\beta(t)}dB_t$

Backward (ODE): $dX(t) = (-2\beta(t)\nabla_X \log P(X, t) - \alpha(t)X(t)) \overline{dt}$

Deterministic Reverse Equation

- If $X(T)$ is initialized with the law resulting from the forward distribution, the marginal of the reverse diffusion are the right ones.
 - No claim on the trajectories. . . but irrelevant in the generative setting.
 - Much faster numerical scheme. . .
-
- Stability results on the score estimation error and the numerical scheme exist for both the stochastic and deterministic case.

$$X \sim P \xrightleftharpoons[P(X|X_1)]{Q(X_1|X)} X_1 \xrightleftharpoons[P(X_1|X_2)]{Q(X_2|X_1)} X_1 \dots \xrightleftharpoons[P(X_t|X_{t+1})]{Q(X_{t+1}|X_t)} \dots X_{T-1} \xrightleftharpoons[P(X_{T-1}|X_T)]{Q(X_T|X_{T-1})} X_T \sim P_T$$

- Generation of X from X_T using $P(X_t|X_{t+1})$.
- Encoder/Forward diffusion: $Q(X_{t+1}|X_t)$.

Variational Auto-Encoder

- P_T is chosen as Gaussian.
- Both generative $P(X_t|X_{t+1})$ and *encoder* $Q(X_{t+1}|X_t)$ have to be learned.

Approximated Diffusion Model

- $Q(X_{t+1}|X_t)$ is known and P_T is approximately Gaussian.
- Generative $P(X_t|X_{t+1})$ has to be learned.
- Denoising tricks can be obtained as an ELBO starting from $Q(X_{t+1}|X_t) = Q(X_{t+1}|X_t, X) \dots$

- 1 Unsupervised Learning?
 - 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - 5 **Generative Modeling**
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - **Generative Adversarial Network**
 - 6 References

$$\omega \sim \tilde{Q}(\cdot|Y) \text{ and } \tilde{X} = G(Y, \omega)$$

Non density based approach

- Can we optimize G without thinking in term of density (or score)?

$$(\bar{X}, Z, Y) = \begin{cases} (X, 1, Y) & \text{with proba } 1/2 \\ (G(Y, \omega), 0, Y) & \text{otherwise} \end{cases}$$

GAN Approach

- Can we guess Z with a discriminator $D(\bar{X}, Y)$?
- No if G is perfect!

$$\begin{aligned} \max_G \min_D \mathbb{E}[\bar{X}, Y] \ell(D(\bar{X}, Y)) \\ = \max_G \min_D \left(\frac{1}{2} \mathbb{E}_{X, Y} [\ell(D(X, Y), 1)] + \frac{1}{2} \mathbb{E}_\omega [\ell(D(G(Y, \omega), Y), 0)] \right) \end{aligned}$$

Discrimination

- Similar idea than the *noise* contrastive approach in EBM.
- If ℓ is a convexification of the $\ell^{0/1}$ loss then the optimal classifier is given by

$$D(\bar{X}|Y) = \begin{cases} 1 & \text{if } p(\bar{X}|Y) > \tilde{p}(\bar{X}|Y) \\ 0 & \text{otherwise.} \end{cases}$$

- If ℓ is the log-likelihood

$$\max_G \min_D \mathbb{E}[\bar{X}, Y] \ell(D(\bar{X}, Y)) = \max_G \log_2 -\mathbb{E}_Y [JKL_{1/2}(p(\cdot|Y), \tilde{p}(\cdot|Y))]$$

- Direct (approximate) optimization using only samples (with the reparametrization trick).

$$\begin{aligned} D_f(P, Q) &= \int f\left(\frac{p(x)}{q(x)}\right) q(x) \\ &= \sup_T \mathbb{E}_{\underline{X} \sim P}[T(\underline{X})] - \mathbb{E}_{G \sim Q}[f^*(T(G))] \end{aligned}$$

f -GAN

- Optimization of

$$\min_G \sup_T (\mathbb{E}_{\underline{X}, Y}[T(\underline{X})] - \mathbb{E}_{\omega, Y}[f^*(T(G(Y, \omega)))])$$

- Direct (approximate) optimization using only samples (with the reparametrization trick).
- Direct extension of the previous scheme.
- T is not a discriminator but there is an explicit link when $f(u) = \log(u)$.

$$\begin{aligned} W(P, Q) &= \inf_{\xi \in \pi(P, Q)} \mathbb{E}_{(p, q) \sim \xi} [\|p - q\|] \\ &= \frac{1}{K} \sup_{\|f\|_L \leq K} \mathbb{E}_{\underline{X} \sim P} [f(\underline{X})] - \mathbb{E}_{G \sim Q} [f(G)] \end{aligned}$$

Wasserstein GAN

- Optimization of

$$\min_G \sup_{\|f\|_L \leq 1} \mathbb{E}_{\underline{X} \sim P} [f(\underline{X})] - \mathbb{E}_Z [f(G(Z))]$$

- Direct (approximate) optimization using only samples (with the reparametrization trick).
- More stability but hard to optimize on all the 1-Lipschitz functions.

- 1 Unsupervised Learning?
 - 2 A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - 3 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - 4 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - 5 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - 6 References



F. Husson, S. Le, and J. Pagès.
Exploratory Multivariate Analysis by Example Using R (2nd ed.)
Chapman and Hall/CRC, 2017



B. Ghogh, M. Crowley, F. Karray, and A. Ghodsi.
Elements of Dimensionality Reduction and Manifold Learning.
Springer, 2023



Ch. Aggarwal and Ch. Reddy.
Data Clustering: Algorithms and Applications.
Chapman and Hall/CRC, 2013



Ch. Hennig, M. Meila, F. Murtagh, and R. Rocci.
Handbook of Cluster Analysis.
Chapman and Hall/CRC, 2015



Ch. Bouveyron, G. Celeux, B. Murphy, and A. Raftery.
Model-Based Clustering and Classification for Data Science.
Cambridge University Press, 2019



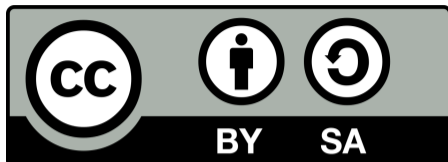
J. Tomczak.
Deep Generative Modeling.
Springer, 2021



D. Foster.
Generative Deep Learning (2nd ed.)
O'Reilly, 2023



A. Géron.
Hands-On Machine Learning with Scikit-Learn, Keras and TensorFlow (3rd ed.)
O'Reilly, 2022



Creative Commons Attribution-ShareAlike (CC BY-SA 4.0)

- You are free to:
 - **Share:** copy and redistribute the material in any medium or format
 - **Adapt:** remix, transform, and build upon the material for any purpose, even commercially.
- Under the following terms:
 - **Attribution:** You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
 - **ShareAlike:** If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.
 - **No additional restrictions:** You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

Contributors

- Main contributor: E. Le Pennec
- Contributors: S. Boucheron, A. Dieuleveut, A.K. Fermin, S. Gadat, S. Gaiffas, A. Guilloux, Ch. Keribin, E. Matzner, M. Sangnier, E. Scornet.