UNSUPERVISED MACHINE LEARNING FOR IMPROVED UNDERSTANDING, OPTIMIZATION, AND PROCESS MONITORING, FAULT DIAGNOSIS IN EARLY DEVELOPMENT

NELSON LEE AFANADOR, PH.D., RICHARD BAUMGARTNER, PH.D., DAI FENG, PH.D., SETH CLARK, PH.D.



Presentation Date: January 30, 2019

Outline

- Background
- Euclid (L2-norm), Manhattan (L1-norm)
- Unsupervised Random Forest (URF)
- Shallow dive into PCA and URF
- Scale-down model comparison with PCA and URF
- Non-linearly separable clusters PCA, Kernel-PCA, URF
- Final thoughts



'Multi-Block' Data Matrix



| t measurements |
|----------------|----------------|----------------|----------------|-------------------|
| Variable 1 | Variable 2 | Variable 3 | | Variable <i>k</i> |



Overarching Goal

Dimensionality reduction via latent variable methods, still rule the day.

Both distance and kernel matrices are ways we can to understand the distribution of batches. 'Latent spaces' help us make this effort possible.

But how do we get here? Turns out there's multiple ways...





The nature of our multiple objectives

- Process/Analytical understanding
- Improvement/optimization
- Root causes for changes
- Process Monitoring

These objectives have been mode more interesting (or complicated) due to automation and the exponential increase in sampling frequencies.



Trends in academic research related to fault diagnosis



"Trends in academic research related to fault diagnosis based on number of publications in the IEEE Xplore digital library from 1991 to 2010" - Aldrich, C., Auret, L., *Unsupervised Process Monitoring and Fault Diagnosis with Machine Learning Methods*, Springer-Verlag London 2013



'Unsupervised' Machine Learning

- Given the explosion of available process data, due to automation, it may work better in the higher-dimensional spaces we will be forced to work in
- It encompasses a wide range of proven methods
- Flexible modeling options via parameter fine-tuning
- Cross-validation (via reconstruction error) is how one chooses an optimal method for implementation
- Can be computationally expensive
- Wherever traditional MVDA is being used, unsupervised machine learning methods can be plugged-in



Common Projection/Dimensionality Reduction Methods

- Principal Component Analysis (PCA) → current most popular Euclidian method
- Unsupervised Random Forest
- Kernel PCA
- Independent Component Analysis
- Locally linear embedding
- t-SNE
- Many others...



Common Projection/Dimensionality Reduction Methods

- Principal Component Analysis (PCA) → current most popular Euclidian method
- Unsupervised Random Forest
- Kernel PCA



What is a latent variable in everyday life?

Happiness

•How do you measure happiness?

- -Again, it is a combination of factors
 - Work-life balance
 - Health (a latent variable itself)
 - Vacation days



What is a latent variable in everyday life?

Happiness

•How do you measure happiness?

- -Again, it is a combination of factors
 - Work-life balance
 - Health (a latent variable itself)
 - Vacation days



Grey shaded area is the 'latent space'

What is the best way to determine the distance between observations? *Turns out there's quite a few ways...*





What is the best way to determine the distance between observations? *Turns out there's quite a few ways...*





'Intuitive' Distance Metrics





























An elegant consequence of using Euclidian distances is that -(double-mean centered)/2 the square of this matrix results in, after mean-centering X, in XX' → PCA → derive loadings



























Geometrically what does it all mean?





Distance Metrics in High Dimensional Space





Examples of a not so 'Intuitive' Distance Metrics



[Unsupervised] Random Forest



[Unsupervised] Random Forest



Repeat this process many times but add in

- 1. bootstrap sample for each tree
- 2. At each node, randomly select a subset of predictors to determine a split
- 3. Grow an un-pruned tree
- 4. Vote across all trees in the forest for classification assignment

In URF, observations that are modeled together in the tree are deemed similar



URF – model/probabilistic based distance



Across *n* bootstrap samples

	x1	x2	x3	x4	x5	Y
Obs i	863	252	951	-301	544	Actuai
Obs 2	— 78 —	81	119	3 55	11	Actual
Obs 3	304	564	247	853	212	Actual
Obs 4	735	197	516	262	115	Actual
Ob s 5	215	813	502	- 587		Actual
Obs 6	137	814	30	842	299	Actual
Obs 7	415	201	355	292	259	Actual

Distance matrix



Biologics Mfg. Simulated Example where we perform a multivariate process scale comparison



'Multi-Block' Data Matrix



| t measurements |
|----------------|----------------|----------------|----------------|-------------------|
| Variable 1 | Variable 2 | Variable 3 | | Variable <i>k</i> |



Simulated Process Trajectory Curves



Data Pre-processing

Principal Component Analysis

- If categorical variables present, perform one-hot encoding
- Scale to unit-variance
- Ready for analysis

Unsupervised Random Forest

- One-hot encoding not necessary for categorical variables
 - Means that categorical variable are assessed as a single contributor to the model
- Ready for analysis



PCA Scores (3 Principal Components)





URF PCoA Scores (3 Principal Coordinates)





Variables Contributing to PCA Clustering



Note differences w.r.t. contribution



Variables Contributing to PCoA Clustering



What did PCA & URF discover?



But wait, there's more & here is where it gets more interesting (i.e. challenging)!



What if I'm dealing with highly non-linearly separable original spaces? Here we bring in Kernel PCA



Kernel PCA



 ϕ :

Why do we care? Because now all we need to do is simply calculate the inner product and *implicitly* map to the higher-dimensional space, i.e., we do not have to map to ϕ and THEN compute the inner product.

\Re^2 \Re^3 $\mapsto (z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$ (x_1, x_2) $k(x_i, x_j) = x_{i1}x_{j1} + x_{i2}x_{j2} + \dots + x_{ik}x_{jk} = XX'$ *Assuming a centered (XX')homogenous kernel Voila, the 'kernel trick'! $k(x_i, x_j)^2 = x_i^2, \sqrt{2}x_i x_j, x_j^2$ svd(XX')

MERCK

Kernel PCA



































Final Thoughts

- Distance metrics abound
- Successful dimensionality reduction is dependent on finding a suitable distance metric in order to carry out the 'optimal' embedding of the original data unto a lower-dimensional space
- The singular use of Euclidian distances (among other L_k-Norms) may be suboptimal
 - scale dependent (like others)
 - curse of dimensionality
 - non-linearly separable structures
 - forces an increase in dimensionality in the presence of categorical variables due to one-hot encoding
- Method in unsupervised machine learning should be explored in parallel with current methods to determine their usefulness



Aldrich, C., Auret, L., *Unsupervised Process Monitoring and Fault Diagnosis with Machine Learning Methods*, Springer-Verlag London 2013

Borg, I., Groenen, P., *Modern Multidimensional Scaling – Theory and Applications*, Springer Science+Business Media New York, 1997

Aggarwal C.C., Hinneburg A., Keim D.A. (2001) On the Surprising Behavior of Distance Metrics in High Dimensional Space. In: Van den Bussche J., Vianu V. (eds) Database Theory — ICDT 2001. ICDT 2001. Lecture Notes in Computer Science, vol 1973. Springer, Berlin, Heidelberg



Nelson Lee Afanador, Thanh Tran, Lionel Blanchet, and Richard Baumgartner^{*} (2016). mvdalab: Multivariate Data Analysis Laboratory. R package version 1.2. Andy Liaw, and Matt Wiener (2002). Classification and Regression by randomForest. R News 2(3), 18--22.

Alexandros Karatzoglou, Alex Smola, Kurt Hornik, Achim Zeileis (2004). kernlab -An S4 Package for Kernel Methods in R. Journal of Statistical Software 11(9), 1-20.

H. Wickham. ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag New York, 2016.



THANK YOU!

