

Master of Applied Data Science

High-level synthesis of machine learning

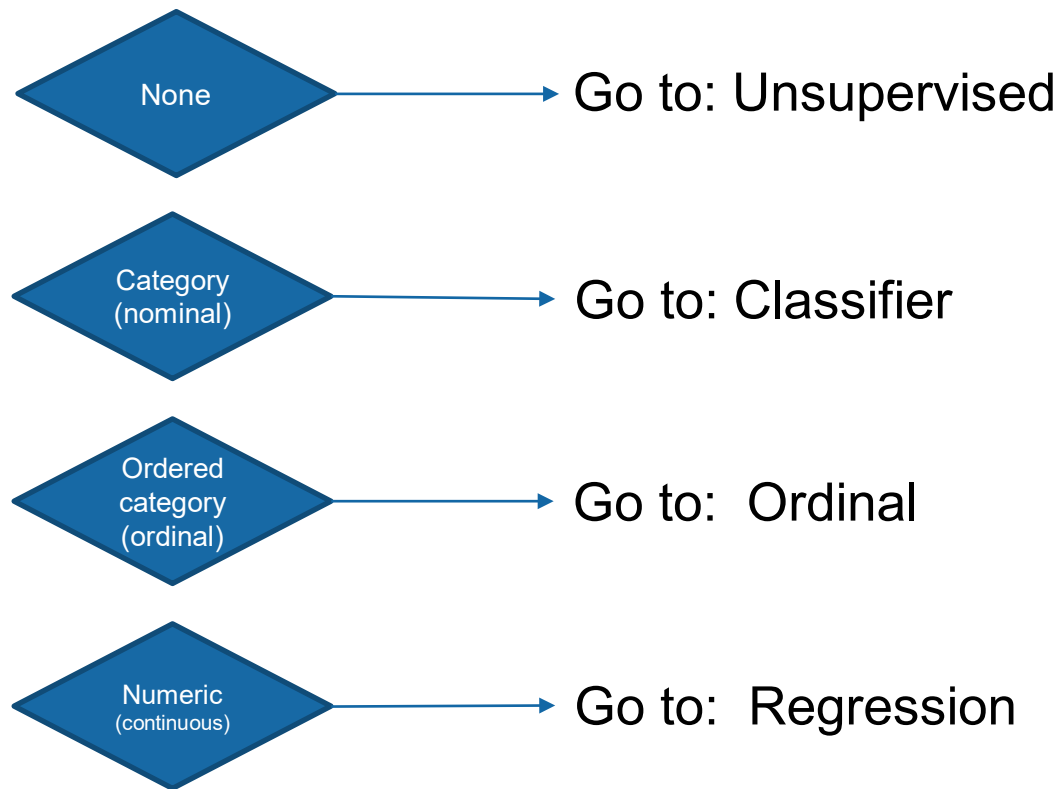
Supervised and Unsupervised Learning
Decision Flowchart:
When to use which method?

Kevyn Collins-Thompson & Yumou Wei

This content not to be redistributed outside of the University of Michigan



What type of target variable
are you predicting?



Start

This version v0.9.2022.08.13

Latest version always at:

https://www.umich.edu/~kevynct/mads_ml_synthesis.pdf

This flowchart is a work in progress.

Send feedback to:

Kevyn Collins-Thompson kevynct@umich.edu

These decision flows are meant as a starting guide only: they should not replace actual thinking about the problem.

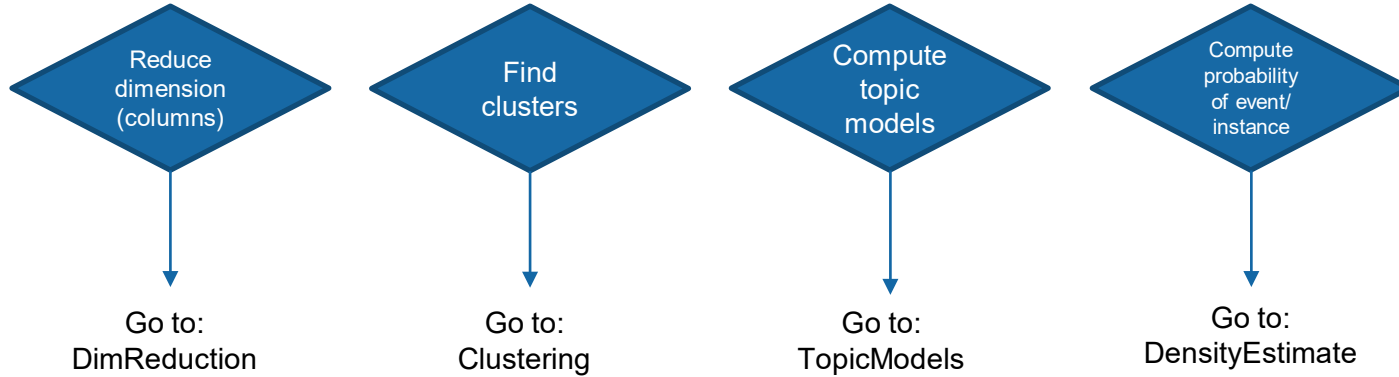
Current assumptions include:

- No missing features or labels
- Roughly balanced classes
- Simple target variable types (e.g. no structures, graphs, permutations, ...)

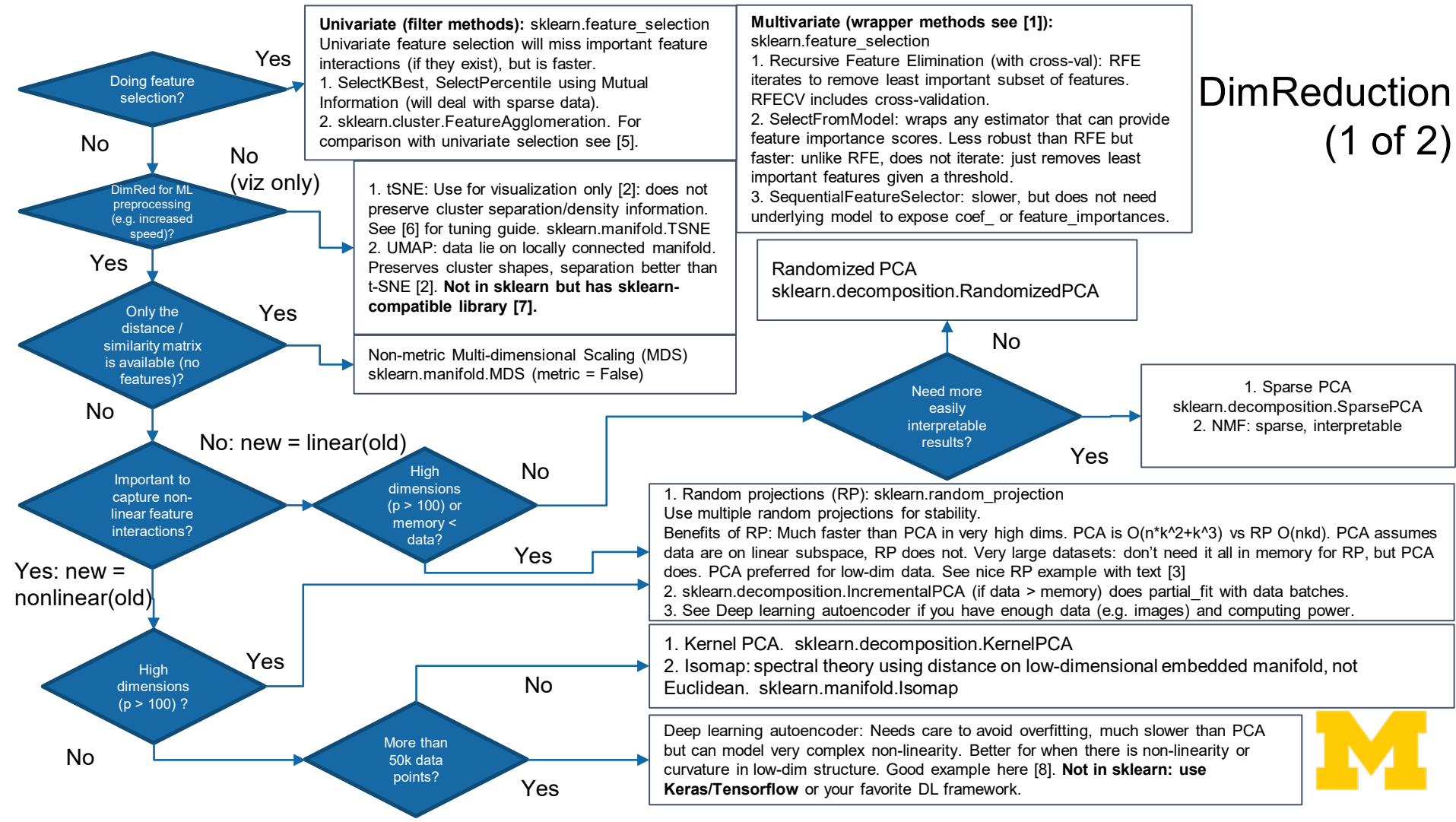


What do you need from your
unlabeled data?

Unsupervised
No target labels



DimReduction (1 of 2)

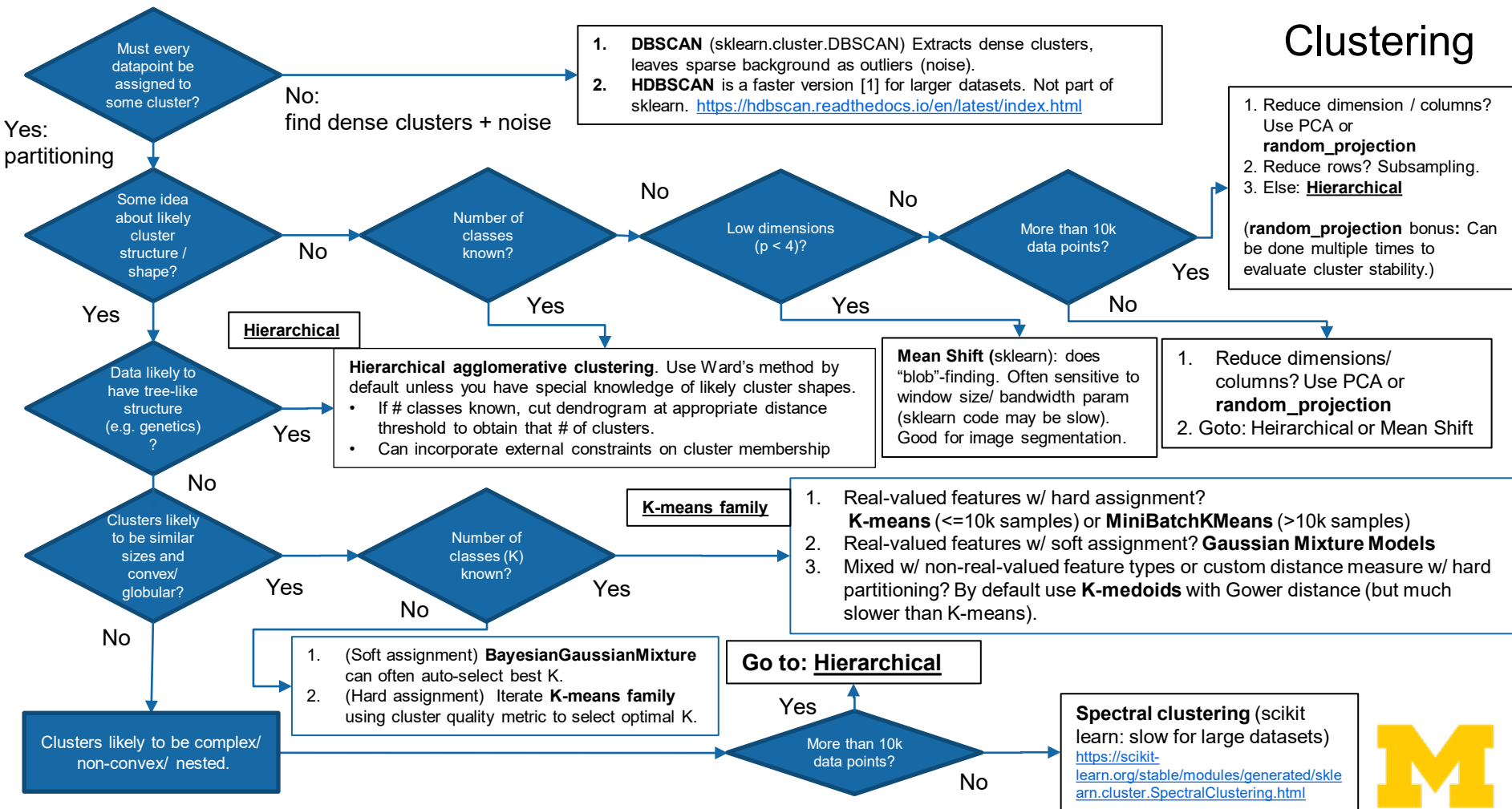


DimReduction (2 of 2)

- [1] R. Kohavi, G.H. John. Wrappers for feature subset selection. <https://ai.stanford.edu/~ronnyk/wrappersPrint.pdf>
- [2] N. Oskolkov. 2020. tSNE vs UMAP: Global Structure
<https://towardsdatascience.com/tsne-vs-umap-global-structure-4d8045acba17>
- [3] B. Schmidt. 2018. "Stable Random Projection: Lightweight, General-Purpose Dimensionality Reduction for Digitized Libraries." *Journal of Cultural Analytics* 3 (1). <https://doi.org/10.22148/16.025>.
- [4] L. Nguyen, S. Holmes. 2019. Ten quick tips for effective dimensionality reduction. <https://doi.org/10.1371/journal.pcbi.1006907>
- [5] https://scikit-learn.org/stable/auto_examples/cluster/plot_feature_agglomeration_vs_univariate_selection.html
- [6] <https://distill.pub/2016/misread-tsne/>
- [7] https://umap-learn.readthedocs.io/en/latest/basic_usage.html
- [8] <https://ekamperi.github.io/machine%20learning/2021/01/21/encoder-decoder-model.html>



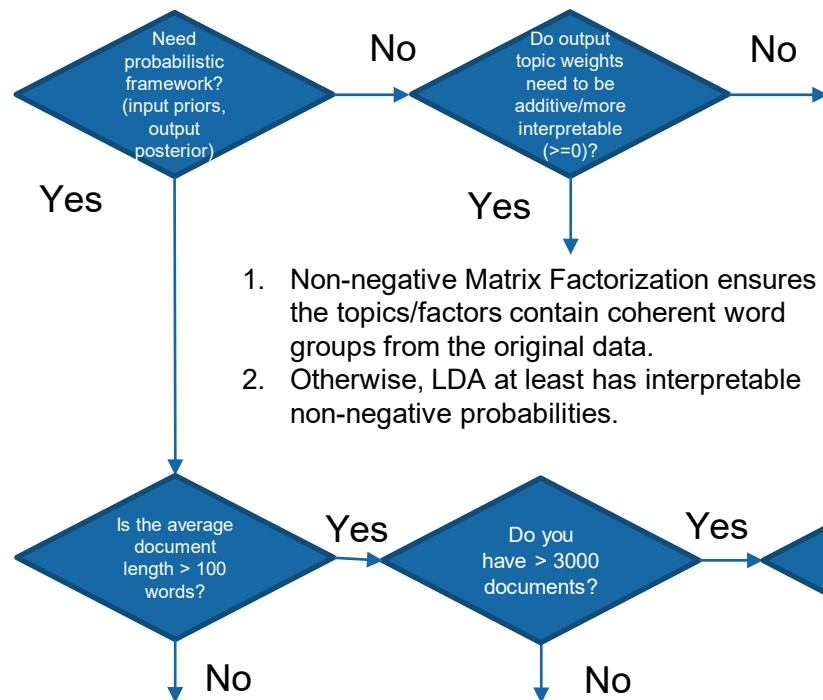
Clustering



[1] https://hdbscan.readthedocs.io/en/latest/comparing_clustering_algorithms.html



Topic Models



Latent Semantic Indexing/Analysis.

- If fit with N topics, can instantly truncate to $K < N$ topics later.
- Term and doc similarity matrices are an easy byproduct, e.g term-term similarity is UU^T , doc-doc similarity is VV^T
- LSA is best at finding a compact semantic representation. LDA best learns descriptive topics [3].

LDA is bad with short docs or too few docs [1].
For short docs consider specialized version called GSDMM (Gibbs Sampling Dirichlet Multinomial Mixture) [2]. <https://github.com/rwalk/gsdmm>

1. Use Latent Dirichlet Allocation, and tune #topics hyperparameter for minimum perplexity and/or maximum coherence.
2. Hierarchical Dirichlet process topic models infer the number of topics from the data. See Gensim: <https://radimrehurek.com/gensim/models/hdpmodel.html>

Source:

[1] Understanding the Limiting Factors of Topic Modeling via Posterior Contraction Analysis. ICML 2014. <http://proceedings.mlr.press/v32/tang14.pdf>

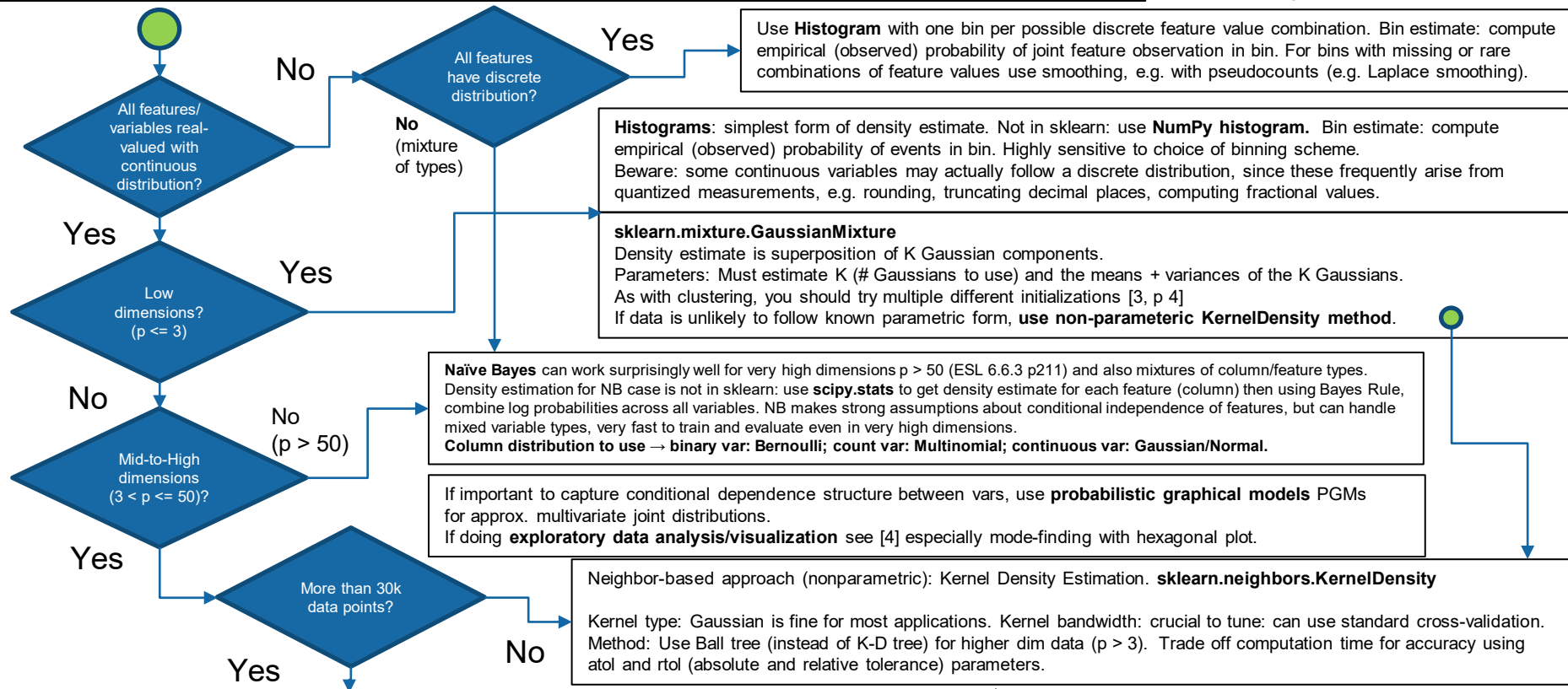
[2] A Dirichlet Multinomial Mixture Model-based Approach for Short Text Clustering. KDD 2014. <https://dl.acm.org/doi/10.1145/2623330.2623715>

[3] Exploring Topic Coherence over Many Models and Many Topics. EMNLP 2012. <https://aclanthology.org/D12-1087.pdf>



Evaluating density estimators: (1) average test negative log-likelihood across cross-validation folds. Lower is better. (2) task-based evaluation based on metric for downstream task that uses the density estimate.

Density Estimation



Naïve Bayes can work surprisingly well for very high dimensions $p > 50$ (ESL 6.6.3 p211) and also mixtures of column/feature types. Density estimation for NB case is not in sklearn: use **scipy.stats** to get density estimate for each feature (column) then using Bayes Rule, combine log probabilities across all variables. NB makes strong assumptions about conditional independence of features, but can handle mixed variable types, very fast to train and evaluate even in very high dimensions.
Column distribution to use → binary var: Bernoulli; count var: Multinomial; continuous var: Gaussian/Normal.

If important to capture conditional dependence structure between vars, use **probabilistic graphical models** PGMs for approx. multivariate joint distributions.
If doing **exploratory data analysis/visualization** see [4] especially mode-finding with hexagonal plot.

Neighbor-based approach (nonparametric): Kernel Density Estimation. **sklearn.neighbors.KernelDensity**

Kernel type: Gaussian is fine for most applications. Kernel bandwidth: crucial to tune: can use standard cross-validation. Method: Use Ball tree (instead of K-D tree) for higher dim data ($p > 3$). Trade off computation time for accuracy using atol and rtol (absolute and relative tolerance) parameters.

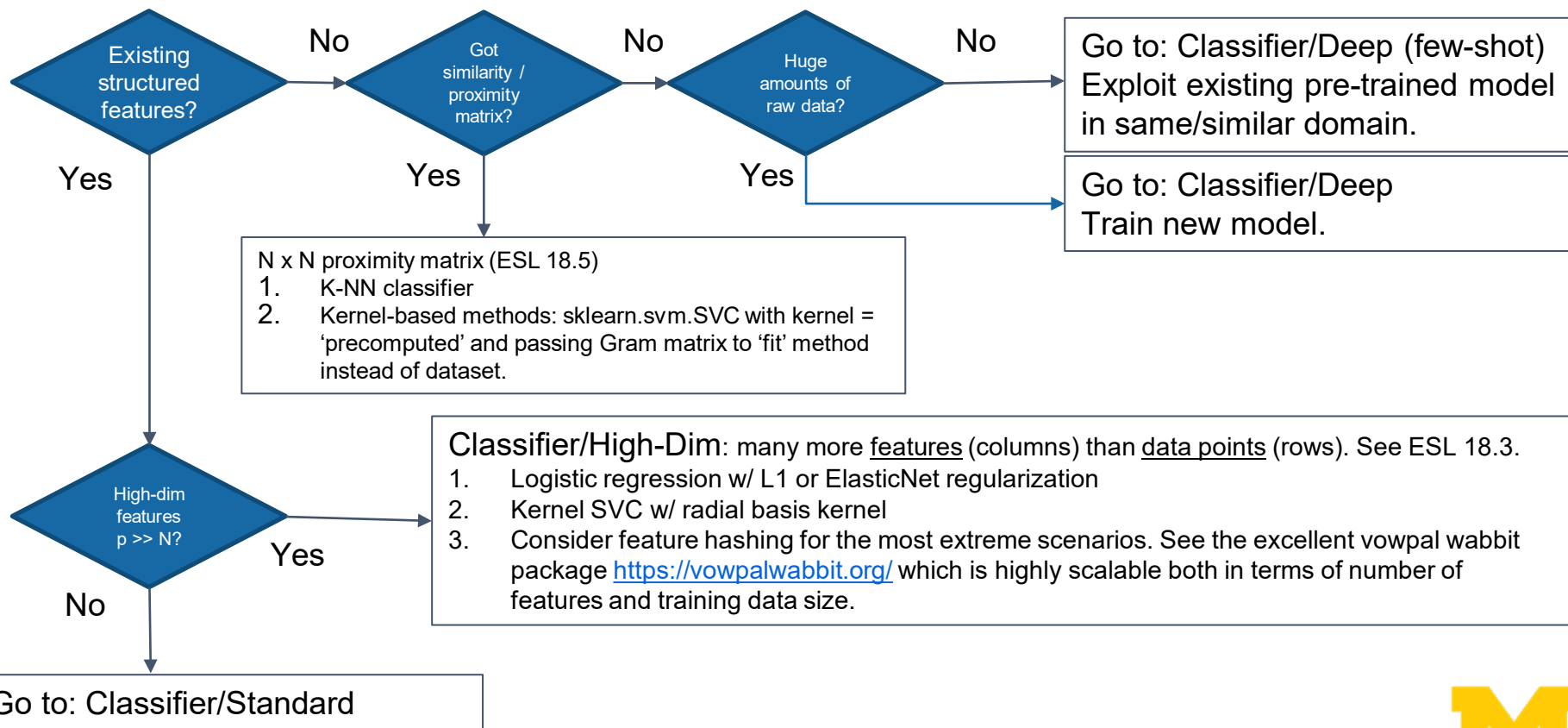
Deep learning: Normalizing Flows (e.g. Real NVP and many other flavors). Transforms a simple density (prior distribution) to a more complex one (posterior distribution) via sequence of invertible transformations [1]. Not available in scikit-learn. One code example using Tensorflow for Real NVP here [2].

ESL refers to Elements of Statistical Learning, 2nd ed. Hastie, Tibshirani, Friedman. Springer.

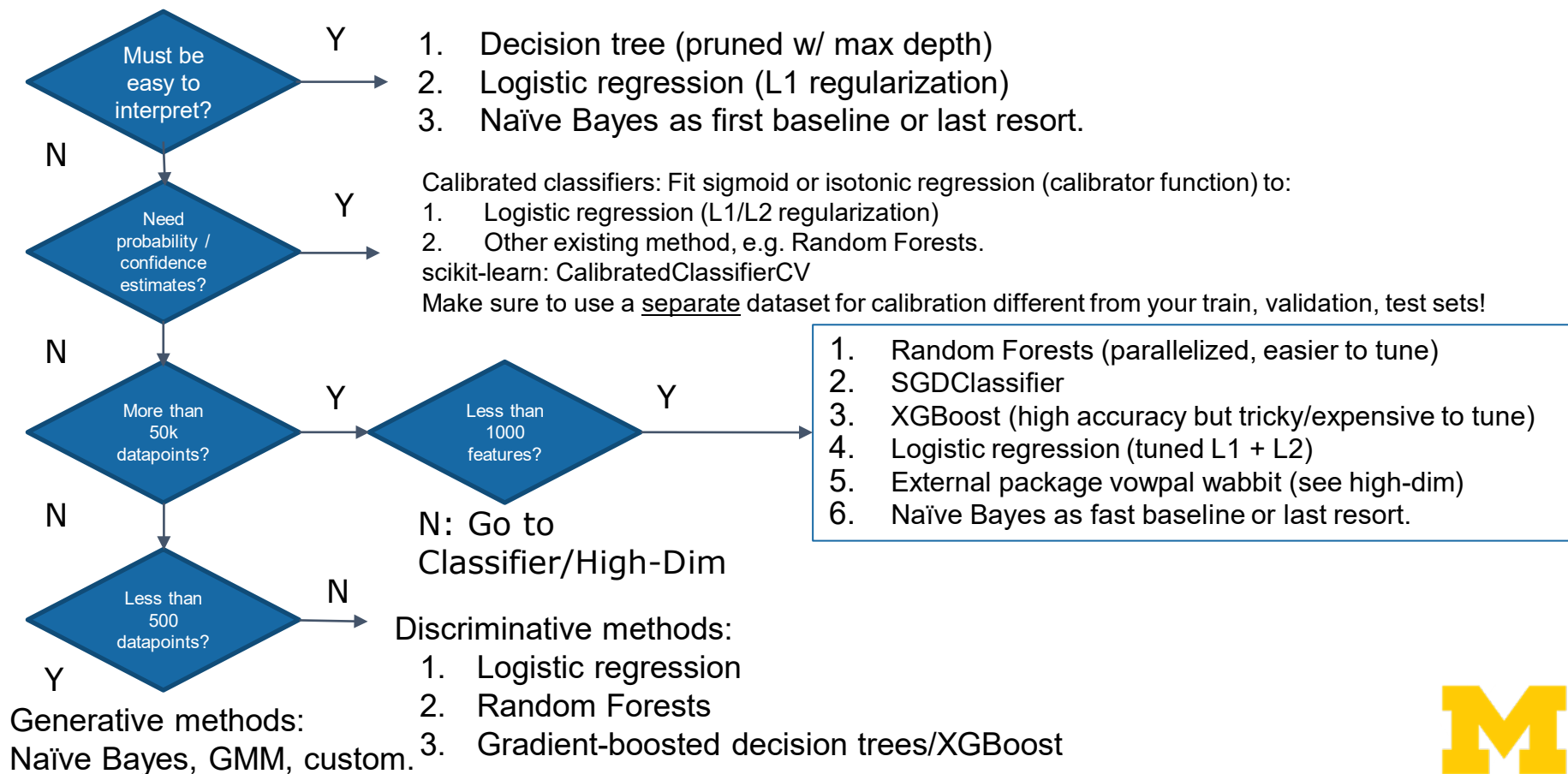
- [1] I. Kobyzev et al. Normalizing Flows: An Intro & Review: <https://arxiv.org/pdf/1908.09257.pdf>
- [2] M. Maria et al. Density estimation using Real NVP. https://keras.io/examples/generative/real_nvp/
- [3] Wang & Scott. Nonparametric Density Estimation for High-dim Data. <https://arxiv.org/pdf/1904.00176.pdf>
- [4] Scott, David W. (2004) : Multivariate Density Estimation and Visualization, Papers, No. 2004,16, Humboldt-Universität zu Berlin, Center for Applied Statistics and Economics (CASE), Berlin https://www.econstor.eu/bitstream/10419/22190/1/16_ds.pdf



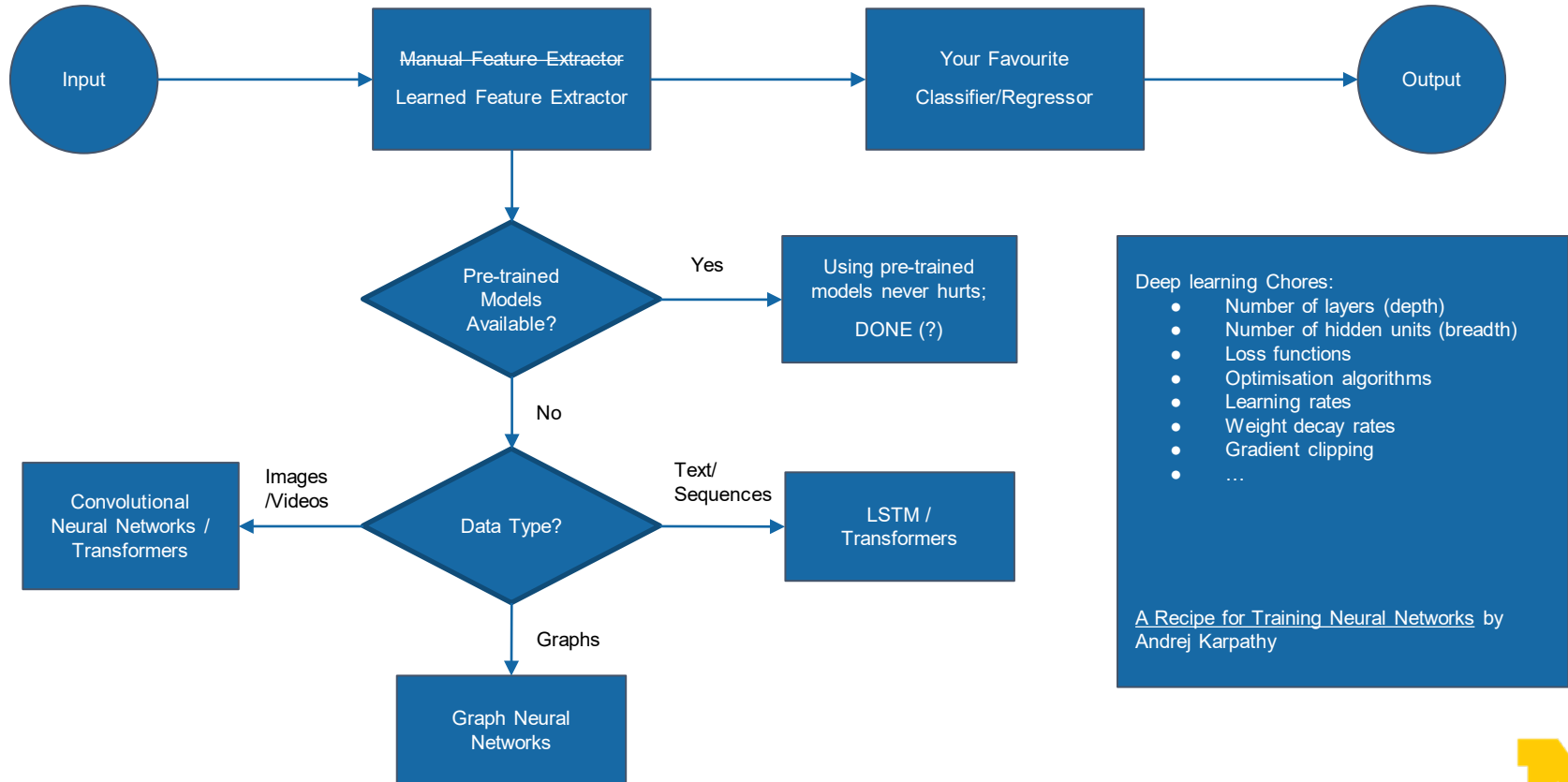
Classifier

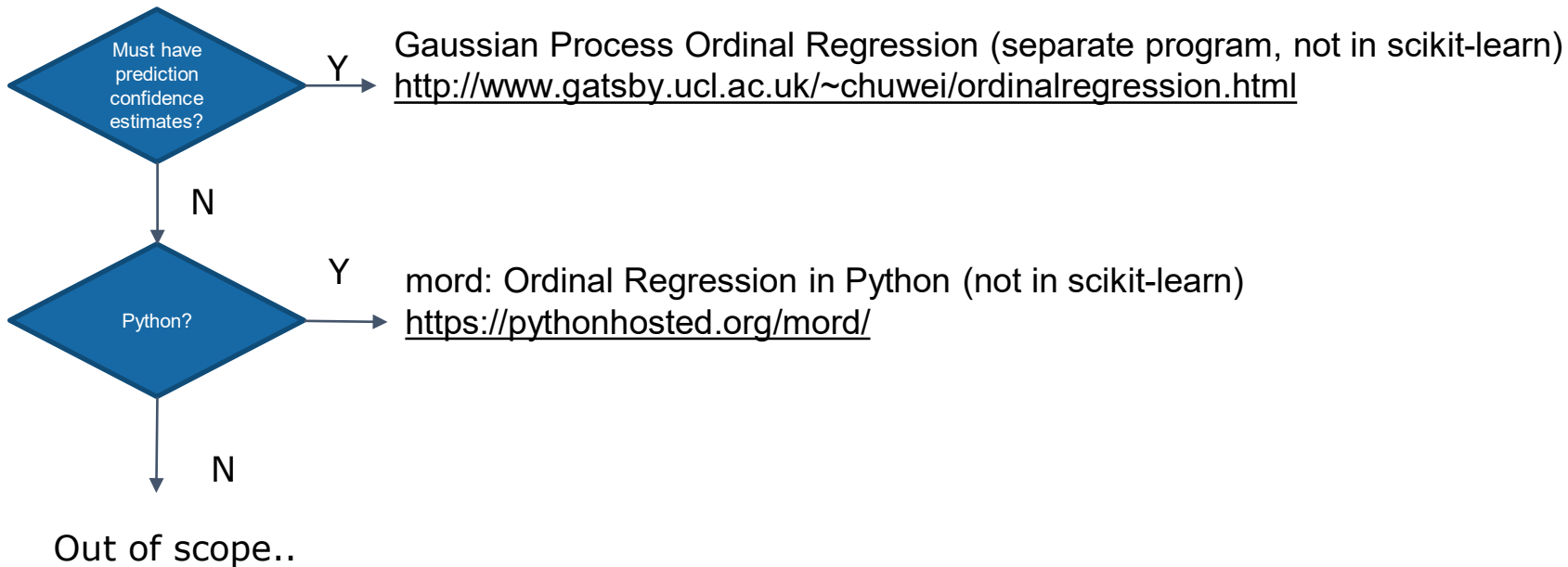


Classifier/Standard



Classifier/Deep





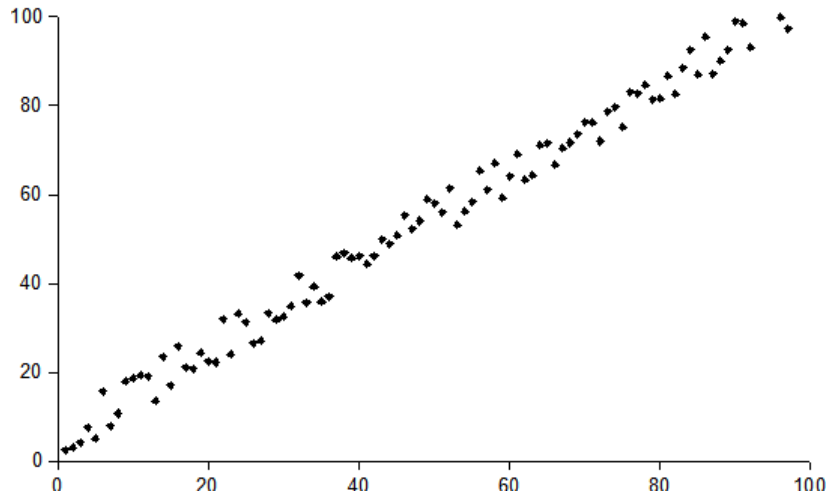
Regression plan

- Background on some regression terms
- How much training data is enough for regression?
- Regression flow chart

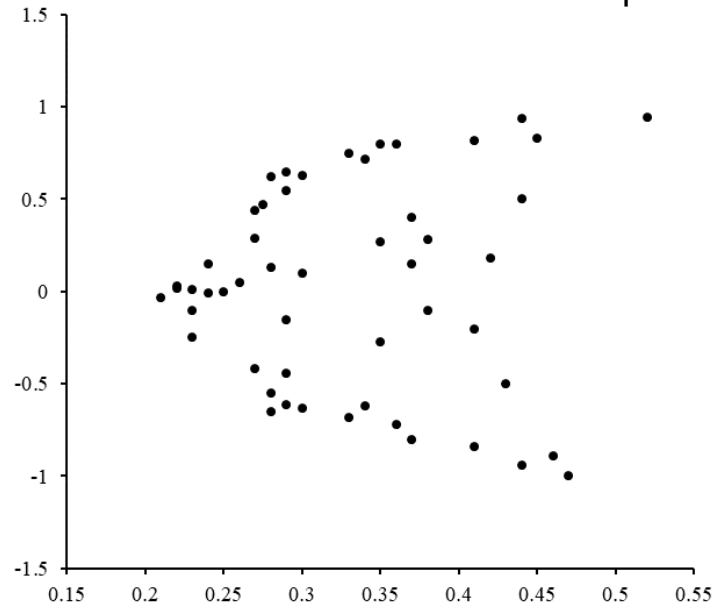


Homoscedasticity (say 3x fast)

Yes homoscedastic: noise is even for all x



Not homoscedastic: noise depends on x

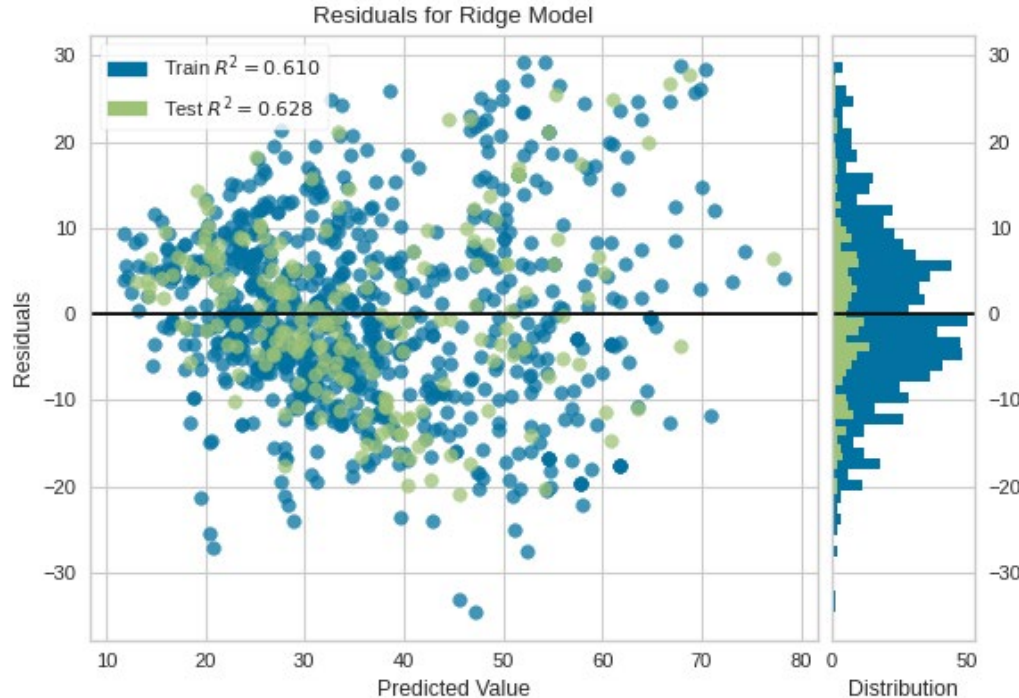


This is one example that motivates the need for enough training data

Source: https://en.wikipedia.org/wiki/Least_squares



Regression residuals



A residual point shows the error in prediction for a single example.

It's the difference between the true (observed) value and the predicted value of the target variable.

A good linear model fit will have

1. Residual points that are randomly dispersed symmetrically around the horizontal axis.
2. A histogram of the above dispersions that is normally distributed around zero.

Source: <https://www.scikit-yb.org/en/latest/api/regressor/residuals.html>



How much training data for regression?

- For ordinary linear regression, 30 rows per parameter/feature is a safe bet.
 - e.g. For a one-independent variable linear regression plus a constant you would need 60 examples.
 - You would only need 30 examples if you know the constant (intercept).
- You need enough points to check assumptions about the data and possibly more robust model selection.
 - Linearity plus homoscedasticity (the variance in noise is same for all observations and doesn't depend on the values of the features)
 - Low correlation of input variables/features (check for multi-collinearity)
 - Checking that residuals are random/normally distributed.
 - Capturing possible interaction effects
 - Doing stepwise variable selection? Double it: 60 examples per parameter/features
- The amount of training data depends on:
 - The specific application
 - How well you know the data, noise and missingness patterns, etc.
 - The objective of the modeling: pure prediction vs confident parameter insights.



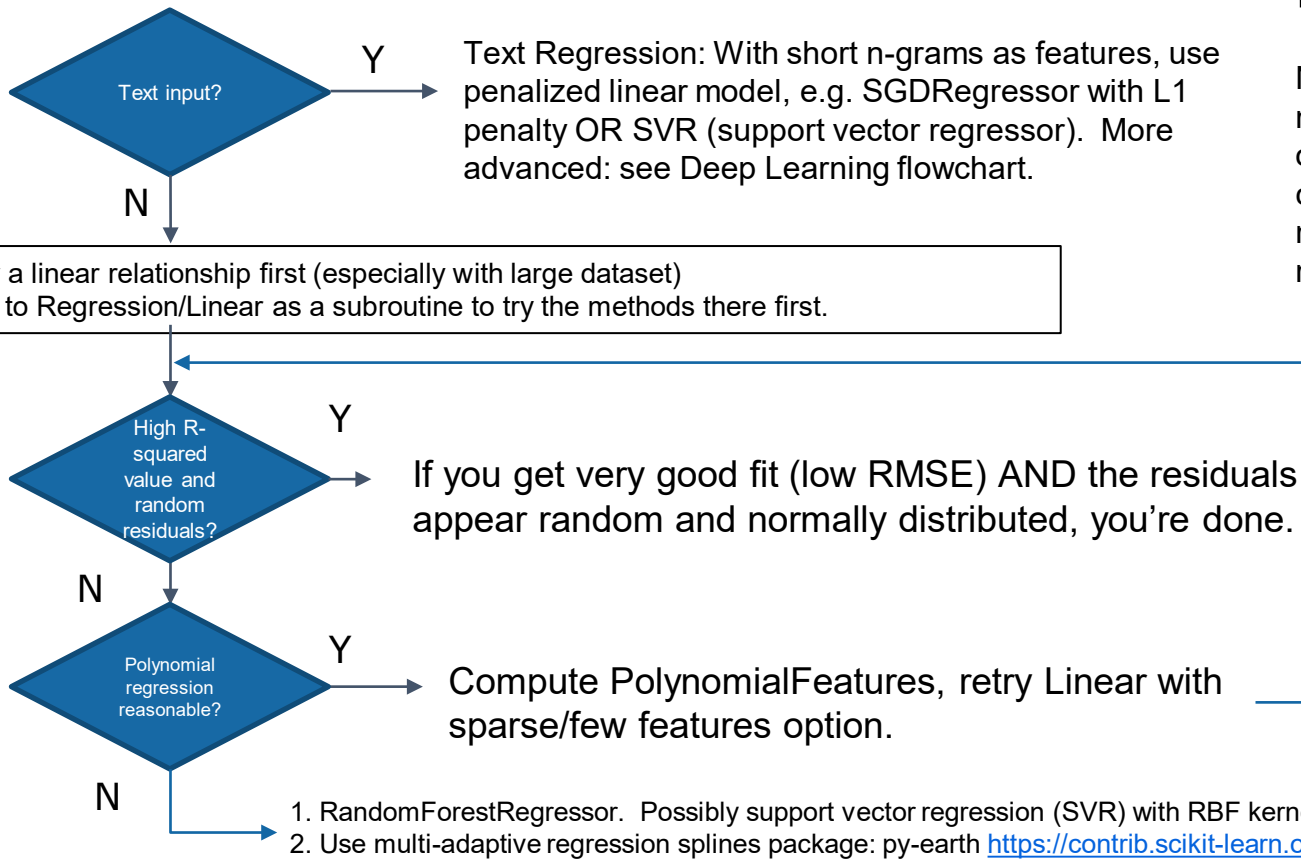
What if I don't have that much data?

- If you have only 10 to 30 examples per feature, use robust alternatives (see chart) unless you have strong theory about the data or some reason to need least squares.
 - HuberRegressor, RANSACRegressor
- With even fewer data points, say only 5-10 per feature you might still use OLS BUT with caveats. You should consider:
 - (a) making the predictions more robust (and reporting confidence intervals!) by using bootstrapping, or use robust regression methods that are less influenced by outliers
 - (b) add additional model assumptions by using Bayesian methods like BayesianRidge with specific priors on the parameters.



Regression

Note: highly simplified decision-making here, and no discussion of hyperparameter tuning or computation constraints. For more in-depth coverage of linear regression see ESL Chapter 2.



Regression/Linear

What are you forecasting about target variable?

1. The median : use mean absolute error (MAE) as the loss metric instead of RMSE.
2. Expected value in some quantile of the target's distribution: use quantile regression.

The mean: use root mean squared error as loss metric

Need probability / confidence estimates?

Y

1. BayesianRidge, especially if you have specific priors on parameters.
2. For general confidence estimates from non-Bayesian regression, use bootstrap estimator with any selected model below (ESL 7.11).

N

Sparse/few features are important?

Y

Use Lasso. Could also use SGDRegressor with either L1 penalty, or ElasticNet penalty, which adds L2 parameter shrinkage like ridge regression, depending on expected nature of sparsity you need and computation speed. To help select the features, you can get the entire Lasso variable selection path using scikit-learn's Lars linear model (least-angle regression: see also ESL 3.4.4)

N

At least 30 examples per feature?

N

At least 15 examples per feature?

N

Y

Ordinary least-squares (OLS):
LinearRegression for uncorrelated features.
SGDRegressor for possible correlated features and large datasets.

Y

Robust least-squares:
e.g. HuberRegression
or RANSACRegression

With only 5-10 examples per feature you could still use OLS BUT with caveats. You should consider (a) making the predictions more robust (and reporting confidence intervals!) by using bootstrapping, or use robust regression methods that are less influenced by outliers; (b) add additional model assumptions by using Bayesian methods like BayesianRidge with specific priors on the parameters.





Kevyn Collins-Thompson

kevynct@umich.edu

School of Information

