

(14) ■ **SVMs:** Maximum margin classifier, decision boundary: $f(x) = \beta_0 + x^T \beta$. Idea is to create a margin on either side of the decision boundary to better separate the classes. ■ **Support Vectors** are the points directly on the margin and inside the boundaries but correctly classified and points that are incorrectly classified. ■ Linear SVMs as empirical risk minimization. Use Lagrange multipliers to reframe $\min_{\beta} \frac{1}{2} \|\beta\|^2 + C \|\beta\|_1$ to $\min_{\beta} \sum_i \text{LOSS} + \lambda \sum_i \|\beta_i\|_1$ called the L0SS + Hinge Loss = $\frac{1}{2} (x - \gamma(\beta + \beta_0))$.

■ **Strength + Weaknesses of Linear SVMs:** (1) Does not depend on Gaussianity or any other probability distribution (2) SVMs tend to more robust than logistic regression such as to outliers as the outliers are only affected by a linear solution whereas in CE the loss would shoot way up and the outliers would have exponential loss. (3) $p \gg n$, SVMs are still computable and have a unique and global solution (b/c regularization). ■ **Non-Linear (Kernel) SVMs:** Recall: kernels and RKHS is take infinite-dimensional non-linear problems and turn them into problems depending on n dimensions (result of Mercer's Theorem). $f(x) = \beta_0 + \sum_i \alpha_i \phi_i(x)$. ■ **Kernel Trick:** Derive optimization problem as a function of inner products $\langle x_i^T x_j \rangle$ in \mathcal{X} , replace all linear inner products with $x_i^T x_j$ with $k(x_i^T x_j)$. ■ **SVM Dual Problem** (equivalent to other forms listed below): $\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j x_i^T x_j y_i y_j$ subject to $\sum_{i=1}^n \alpha_i y_i = 1, 0 \leq \alpha_i \leq C, i = 1, \dots, n$. Param $\alpha \in \mathbb{R}^n$. Hyperparameter $C \geq 0$. The kernel SVM is just taking $x_i^T x_j$ and replacing with $k(x_i^T x_j, x_j)$, everytime we use an inner product to replace it with the kernel of the inner product. ■ **Interpret:** α are sparse and $\alpha_i > 0$ are support vectors. Linear: $\beta = \sum_{i=1}^n \alpha_i x_i^T \beta_0$. Kernel: Predictions $x^T \rightarrow f(x) = \sum_{i=1}^n \alpha_i y_i k(x, x_i)$, predictions are only a function of the support vectors. ■ **Kernels:** RKBF kernels are more flexible than polynomial kernels.

(16) ■ **Unsupervised Learning:** Only have data X_{data} , without y labels/outcomes. Goal: find some structure in the data that is likely to generalize to future data. PATTERNS via dimension reduction techniques, GROUPS via clustering techniques, ASSOCIATION via graphs/networks, ANOMALIES, ...

■ **Dimension Reduction:** Have p features and want to reduce these/ engineer new features $r < p$ that retain major features of the data. Q: Why use this? (1) To visualize the data o/w. w/o dimension reduction we have $\binom{p}{2}$ scatter plots? (2) Want maybe because there is lots of redundancy or irrelevant features (denoising). (3) Pattern discovery. Note that this is useful for both supervised AND unsupervised learning: However, pattern discovery leads towards unsupervised learning.

■ **Principal Component Analysis (PCA):** Note that PCA is the optimal linear reduction technique! $X = UZV^T$ where ZV^T gives principal component scores (e.g. embedding of chapters) and V gives the direction. Idea is to retain high variance patterns via linear projections. Linear combinations of features that maximize variance. Orthogonal projection of points onto hyperplane such that the variance is maximized. These high variance patterns retain the most information in lower dimensions! Q: Is the PC line the same as LS line? No, the OLS takes the vertical projection, not orthogonal projection. Set up PCA mathematically.

■ **Optimize:** Line/hyperplane = $Xv, v \in \mathbb{R}^p$ (weights on each feature) and if it is orthogonal, $v^T v = 1$. Then Xv is a linear projection and Xv becomes the first PC line/plane. Now we want to find the weights v that maximizes the variance of this plane. $\max_v \text{Var}(Xv)$ subject to $v^T v = 1$. Then $\max_v \text{Var}(Xv) = \max_v \text{Var}(X^T X v) = \max_v v^T X^T X v$ where $X^T X$ is the variance-covariance matrix; we don't actually know the population $X^T X$ so instead we will use an estimate of $X^T X$ (assuming $X \sim N(0, \Sigma)$) then $Z = \frac{1}{\sqrt{n}} X^T X$ (MLE of Σ). Therefore, the optimization problem is not $\max_v v^T X^T X v$ subject to $v^T v = 1$ (to find the first PC) and subsequent PC must be orthogonal. All PC: $\max_v v^T X^T X v$ subject to $v_1^T v_j = 1$ and $v_i^T v_j = 0 \forall k \neq j$ and $k = 1, \dots, r$. These constraints are to ensure orthogonality.

■ **Solution:** Eigen decomposition. Recall, $v_1^T X^T X v_1 = \lambda_1$. PC solution is given by Eigen Decomposition of $X^T X$. Alternatively given by Singular Value Decomposition of $X = SVD(X)$. ■ **SVD:** $\max_v v^T X^T X v = \max_{\alpha, u} u^T u = 1, v_1^T v_1 = 0, u_1^T u_1 = 0, \forall k \neq 1$. Review SVD, $X_{n \times p} = U_{n \times n} \Sigma_{n \times p} V_{p \times p}^T$, U, V are orthonormal ($U^T U = I$) UNSEALD PC, $V_{p \times p}$ are orthonormal PC DIRECTIONS, $D_{p \times p}$ diagonal matrix $d_i \geq d_2 \geq \dots \geq d_p$, $\text{rank}(X) = \text{PROF of THE VARIANCE of EACH PC}$. How to interpret SVD solution: v_1 is the first major observation pattern and v_1 is the first major feature pattern. u_1 are subsequent patterns are uncorrelated and note that d_i is the first singular value and the variance explained by the first pattern is d_i^2 . ■ **PVE (Proportion of Variance Explained):** $\frac{d_i^2}{\sum_{j=1}^p d_j^2}$. ■ **CVE (Cumulative VIE - screeplot):** $\frac{\sum_{j=1}^i d_j^2}{\sum_{j=1}^p d_j^2}$ by the first k PCs. ■ **Properties of PCA:** (1) The PC problem is NON-CONVEX, yet there exists a global unique closed form solution solved by SVD!!!!!!!!!!!!!!!!!!!!!!!!!!!!!! (2) u, v are unique up to a sign change. (3) The solution is ordered and nested (important for interpretation), i.e., solution is the same regardless of number of PC selected, r dimensions.

(17) ■ **PCA (Cont.):** PCA is the best linear dimension reduction technique, $\min_{X'} \|X - X'\|_F^2$ subject to $\text{rank}(Z) = k$ yields the SVD solution! PCA is very good for interpretation. ■ **Practicalities:** To center and scale or not? Application dependent 2-models: (Application 1) - Cov model $X \sim N(\mu, \Sigma) \rightarrow \max \text{Var}$. This model implies we should center the features! If we scale the features, then we maximize variance patterns in the correlation matrix instead of variance-covariance matrix. If the scales are completely different such as income and proportion of race or something, then scaling could make sense! (Application 2) - Low-Rank Mean Model, $X = LR + E$ where $E \sim F(\mu, \sigma^2)$. We don't want to center or scale here, our goal is to find the mean! Inference: If we don't assume multivariate Gaussianity then we might not want to center or scale our data. It is okay to try with and without centering and choose which has better patterns but don't scale unless the scales of features differ significantly.

■ **PCA Extensions:** Regularized PCA \rightarrow functional PCA, PCA for smooth factors, spatial + timeseries. Sparse PCA \rightarrow L_1 penalty. Supervised PCA (such as LDA), etc.

■ **Other Linear Dimensional Reduction Techniques:** (Although PCA is optimal linear DR, sometimes we may not want to find the highest variance patterns). We may want to interpret different types of patterns beyond PCA. ■ **Non-Positive Matrix Factorization (NMF):** $\min_{X'} \|X - X'\|_F^2$ subject to $L_{ij} \geq 0$ and $R_{ij} \geq 0$ and $X_{ij} \geq 0$ so this is exactly PCA except instead of orthogonality constraints we are adding non-negativity constraints! Why do we want to do this? NMF finds patterns that go in one direction, i.e., PCA finds orthogonal directions that maximize variance, but components can be positive and negative features, making interpretation harder. NMF finds additive, non-negative directions (parts-based patterns), making it easier to interpret components as intuitive groups or traits. These are called **Archetypal patterns**, not optimal mathematically but very interpretable. ■ **Theory:** NMF yields local solution (you can get different solutions from different starting patterns) and is not nested/ordering - meaning choosing k is super important (change k yields diff. solts). This basically yields a set of patterns that seem significant but cannot determine which is most important/significant.

■ **Non-Linear Dimension Reduction (Manifold Learning):** Why do we want nonlinear DR? The most interesting patterns don't necessarily have to lie on a hyperplane! Project our data onto a manifold - some geometric object in r -dimensional space. Then we can visualize the data (n observations) as projected onto the manifold by unfolding this manifold into a hyperplane! ■ **Manifold learning** seeks a lower dimensional representation $X' \in \mathbb{R}^r, r \ll p$, that represents the patterns/differences amongst observations.

■ **Limitations:** In Manifold learning, (1) we only get observation patterns and we cannot get/interpret feature patterns unlike Linear DR (2) We never learn the manifold. ■ **Manifold Map:** This is the visualization of the manifold as a hyperplane. These are very good for interpretation. They are invariant to rotation or translations meaning the position does not matter because we can unfold or flip the manifold before we use a hyperplane - we only care about the relative distances between observations in manifold map. NOTE: There are two types of manifold learning (1) Global - globally represent distances or dissimilarity in lower-dimensional space; can interpret distances between two points as true distance, i.e., if two points are far we can interpret them as far! (2) Local - only locally represent distances or dissimilarities in lower-dimensional space; tries to preserve nearest neighbors so within the neighborhood, this cluster of points are close but the distance between clusters cannot be interpreted as either close or far!

■ **Kernel PCA (GLOBAL):** $D_{n \times n}$ distance or dissimilarities, $D_{ij} = k(x_i, x_j)$ (rbf, poly, etc). Kernel PCA is the eigen decomposition of D (global soln) \rightarrow non-linearize distances and apply PCA.

(18) ■ **Spectral Methods:** Spectral methods use the eigenvectors of similarity or distance matrices to embed data in lower dimensions, capturing global structure (e.g., in PCA, spectral clustering, or Laplacian eigenmaps). **Steps:** (1) Compute $D_{n \times n}$ distance or dissimilarity matrix between observations. (2) Do PCA on $D_{n \times n} = UDU^T$ and the first k columns of U are the Z Distances. ■ **Distance:** $D_{ij} = \|x_i - x_j\|_2$ are centered distances. How is Classical MDS similar to PCA? ■ **Theorem:** Classical MDS is exactly equivalent to PCA for Euclidean distances. Note, you can also use other distances beyond Euclidean distances such as L_1 , Hamming, ...

■ **Spectral Methods:** PCA, Kernel PCA, Classical MDS, Spectral Embedding, Laplacian Eigenmaps, Spectral Clustering. **Not spectral methods:** t-SNE, UMAP, MRF. ■ **Spectral Embedding:** D as a graph-based similarity matrix (turn our observations into a $n \times n$ graph, if two observations are close draw edge, threshold = epsilon graph, nearest neighbor graphs - connect closest neighbors). $A_{n \times n}$ is the adjacency matrix (symmetry, so almost want inverse of this) and $A_{ij} =$ edge between observations i and j . Use Normalized Laplacian (like inverse of $A_{n \times n}$ acts like dissimilarity for graphs) $L_{n \times n} = D_{n \times n} - A_{n \times n}$ where $D_{n \times n}$ is diagonal matrix of degrees for each node.

■ **Global (Non-Spectral) Manifold Methods:** Metric MDS, seeks to optimize a loss function that keeps distances in X (original data) and Z (dimension reduced data) close! Therefore, $\min_{Z'} \sum_{i,j} \|x_i - x_j\|_2 - \|z_i - z_j\|_2^2$ (Euclidean distance in lower dimensional space + clearly tries to preserve global distances). CONVEX global solution but we DO NOT get a nested and ordered solution.

■ **Local (Neighbor Embedding) Methods:** Find a lower dimensional space ($Z_{n \times k}$) such that close neighbors remain close, i.e., not all distance are preserved just within cluster distances are preserved. Only interpret locally! ■ **t-Stochastic Neighbor Embedding (t-SNE), Uniform Manifold Approximation Projection (UMAP)**

■ **Steps:** (1) Calculate the original space embedding probabilities (normalized dissimilarities), e.g. for t-SNE: $p_{ij} = \frac{e^{-\frac{\|x_i - x_j\|_2^2}{2\sigma^2}}}{\sum_{j=1}^n e^{-\frac{\|x_i - x_j\|_2^2}{2\sigma^2}}}$ and UMAP uses generalized distances and calculates distances only over nearest neighbors. (2) Define lower dimensional space embedding probabilities $(q_{ij}, z_{ij}) = g(\|z_i - z_j\|_2)$ where t-SNE \rightarrow Cauchy and UMAP \rightarrow generalized Cauchy. Far away points get pushed further away! (3) Learn Z by fitting some loss function between original and lower dimensional probabilities. t-SNE uses KL-divergence and UMAP uses Cross Entropy Loss $= \min_{Z'} \sum_{i,j} p_{ij} \log(\frac{q_{ij}}{p_{ij}}) + (1 - p_{ij}) \log(1 - \frac{q_{ij}}{p_{ij}})$.

■ **PreCons:** Very useful for visualizing local patterns and finding clusters. Although they yield a local solution and are very sensitive to hyperparameter selection.

■ **Application:** Interpret UMAP with caution and JOINTLY USE THIS METHOD WITH A GLOBAL METHOD like PCA!!!!!!

■ **Clustering:** Want to find groups of observations (or features or both) that are similar. ■ **Flat vs Nested Clusters:** A nested cluster is that smaller groups are preserved within the larger group, i.e., if we changed the hyperparameter for number of clusters it would just merge groups to create bigger clusters. Flat is just not-nested; note that flat clustering is very dependent on the hyperparameter K (the number of clusters)! ■ **Notation:** K is the number of clusters (fixed) and $C(1, \dots, n) = \{1, \dots, K\}$ where $C(i)$ gives the cluster assignment for obs i and d_{ij} is the distance lower observation i and j . ■ **Goal/Optimization:** Min the within cluster dist $\rightarrow \min_{C'} \sum_{i,j \in C} \sum_{i,j \notin C} d_{ij} \cdot \mathbb{1}(C(i) \neq C(j))$. Group similar obs. This is NON-CONVEX (np-hard and computationally infeasible for medium sized data). Therefore, we need to find local approximations for clustering!

(19) ■ **Kmeans Clustering:** Gives hard clustering labels; assumes K . Goal is to approximately solve $\min_{C'} \sum_{i,j \in C} \sum_{i,j \notin C} d_{ij} = \min_{C'} \sum_{i,j \in C} \sum_{i,j \notin C} \|x_i - x_j\|_2$ in a fast computational manner. Use Euclidean distance (Kmeans only uses Euclidean distance). Simplify this optimization by using the mean of cluster $k \rightarrow \bar{x}_k = \frac{1}{n_k} \sum_{i \in C(k)} x_i$ into $\min_{C'} \sum_{i,j \in C} \sum_{i,j \notin C} \|x_i - \bar{x}_k\|_2$ (still NP-hard). Now we apply the **Kmeans trick:** Auxiliary variable, mean of cluster k $\hat{x}_k \rightarrow \min_{C'} \sum_{i,j \in C} \sum_{i,j \notin C} \|x_i - \hat{x}_k\|_2$ (EQN 8) ■ **Kmeans Algorithm** (iterative algo): Pseudo code: some initialization of \hat{x}_k . While not converged: (a) minimize (EQN 8) over \hat{x}_k and minimize (EQN 9) over m_k given $C()$, and (b) converged. ■ **Mathematically:** (b) $\min_{C'} \sum_{i,j \in C} \sum_{i,j \notin C} \|x_i - m_k\|_2 = \hat{x}_k$ (sample mean of the k th cluster). (a) $\min_{\hat{x}_k} \sum_{i,j \in C} \sum_{i,j \notin C} \|x_i - \hat{x}_k\|_2$ is harder then (b), but still simple, it is simply applying Nearest Centroid Classifier (assign points that optimize the assignment - minimize Euclidean distance - to its nearest mean m_k). Geometric solution, e.g., for 3 classes you draw a triangle with vertices at the means and edges connecting the vertices and take a decision boundary orthogonal to the midpoint of each edge. ■ **Summary:** Kmeans algorithm iterates between taking a sample mean and Nearest Centroid Classifier \rightarrow very simple and fast; good for big data.

■ **Properties:** (1) Attain local solution (convergence guaranteed when means stabilize), can't certify a global solution. (2) Linear cluster solution! (3) Very dependent on initialization, therefore we can use Kmeans++ initialization spreads out initial centroid. Worst initializations happen when all the centroids are close together and the clusters don't separate out well. (4) Kmeans will perform best when NCC is optimal or equivalently Naïve Bayes Classifier under certain conditions, i.e., Kmeans performs well under Gaussian + spherical covariance (Kmeans finds clusters that are spheres/uncorrelated features and tends to find balanced clusters (of equal sizes). Dependent on K (not nested); changing K will often dramatically change the solution. ■ **Pro/Cons:** Advantages are quick and intuitive. Cons are the restrictive assumptions and sensitivity to K ; choosing K via validation is super important! ■ **Curse:** When $p \gg n$ there exists the curse of dimensionality makes calculating distances very difficult therefore Kmeans will perform suboptimal if directly applied - we must first apply a dimension reduction such as PCA (to spread out the data into clusters and sphere the data by decorrelating). Therefore, in high dimensional settings or with lots of correlated features PCA and Kmeans combination is very strong for clustering; assumption check: viz in 2d-PC1, PC2.

■ **Mixture Models:** Soft cluster labels using a generative clustering distribution (assumes underlying distribution of the data given cluster k). ■ **Gaussian Mixture Models:** This is a soft cluster extension of Kmeans! Assume the data comes from a mixture of k different distributions. Assume x_i is the cluster label for the i th observation, $x_i \sim P(x_i = k) = \text{mixture probability } \varphi_2^k(x_i = 1)$. GMM assumes that $x_i | x_i = k \sim N(\mu_k, \sigma_k^2)$ (assume here that we have a spherical covariance). Estimation via MLE! $p(x) = \prod_{i=1}^n p(x_i) = \prod_{i=1}^n \sum_{k=1}^K \varphi_2^k(x_i) \cdot \varphi_2^k(x_i - \mu_k)$. Hence, $\{(\mu_k, \sigma_k^2)\}_{k=1}^K \rightarrow \Sigma$, then $\log p(x) = \sum_{i=1}^n \log \varphi_2^k(x_i - \mu_k) - \log \varphi_2^k(x_i - \mu_k)$ - multiplied by two parameters that we are trying to optimize. ■ **Trick:** It is easy to compute the MLE if Z_i is known, just standard computing. So let's assume that Z_i is a latent or hidden variable \rightarrow auxiliary variable \rightarrow ■ **Expectation-Maximization (EM) Algorithm** (iterative algo): (1) E-step \rightarrow infer the hidden parameters given the parameters, $\pi_i | \bar{x}_i, \hat{\sigma}_i^2$

(2) M-step \rightarrow estimate the parameters given the inferred hidden variables $\pi_i, \mu_i, \sigma_i^2 | \bar{x}_i$. ■ **E-step:** $p(x_i = k | \bar{x}_i, [\bar{x}_i, \mu_i, \hat{\sigma}_i^2]) \rightarrow \delta_{ik} = \frac{\pi_i \varphi_2^k(\bar{x}_i - \mu_i, \hat{\sigma}_i^2)}{\sum_{j=1}^K \pi_j \varphi_2^j(\bar{x}_i - \mu_j, \hat{\sigma}_i^2)}$, hence we get $\hat{\delta}_k$ which is called the responsibilities (same for all mixture models but not GMM) that cluster k takes for observation i which is just the soft cluster assignment for observation i . ■ **M-step:** Is the weighted MLE of the Gaussians: $\bar{x}_k = \frac{1}{n_k} \sum_{i=1}^n \delta_{ik} \cdot \bar{x}_i$ (weighted sample mean), $\hat{\sigma}_k^2 = \frac{1}{n_k} \sum_{i=1}^n \delta_{ik} \cdot \|\bar{x}_i - \bar{x}_k\|_2^2$ (weighted sample variance).

■ **Kmeans and GMM are Linear Clustering Methods!** These methods require the use of Euclidean Distance!!!!!! Cannot use other distance metrics.

(20) ■ **Properties of EM Algorithm:** (1) Monotonically increases the observed data log-likelihood $\ln(p(x_i | \mu_k, \sigma_k^2))$ meaning it always converges! (2) Gives local solution, i.e., it depends on the initialization. For GMM, we initialize to the Kmeans solution. (3) Converges slowly; requires many iterations. (4) Depends heavily on K . ■ **Q:** What if the clusters have non-linear decision boundaries? A: Perform Kmeans after non-linear decision reduction! ■ **Spectral Clustering:** This is just applying Kmeans after applying Spectral Embedding as a dimension reduction. Kmeans to the k smallest eigenvectors of L that don't correspond to all 1×1 eigenvector \rightarrow If a graph has k connected components, then the k smallest eigenvectors of L (other than 1×1) give the indicators of the components.

■ **Properties:** (1) Very nonlinear clusters. (2) Flexible. (3) Works for different shaped clusters. ■ **Hierarchical Clustering:** NESTED CLUSTERING. Family of clusters $K = 1, \dots, K, n$. Visualize nested clusters using a dendrogram - binary rooted tree. Shows family from $K = 1$ at the top of the tree to $K = n$ at the bottom of the tree; terminal leaf - has a single sample in each cluster. Building a dendrogram GREEDY ALGORITHM. Top down clustering (starting at bottom of tree and moving up), bottom up clustering (starting at top of tree and moving down). The order of merging on dendrogram is represented by the height of the merge in the tree! There are $n - 1$ steps/levels on dendrogram. ■ **Properties:** (1) such as GREEDY (approx.) \rightarrow local solution. (2) Height matters \rightarrow distance or dissimilarity at which objects are merged. (3) X-axis has ordering and matters but can be flipped/reflected. (4) The number of clusters are given by a horizontal "cut" across the dendrogram - good interpretation of cluster assignments! ■ **Details:** Input, $D_{n \times n}$ distance/dissimilarity matrix. Q: How to merge objects? A: Use LINKAGES \rightarrow dissimilarity metric used to merge objects of one or more observations. ■ **Linkage:** How we merge sets of observations! ■ **Single Linkage** (Min), $d(A, B) = \min_{i \in A, j \in B} d_{ij}$. ■ **Complete Linkage** (Max), $d(A, B) = \max_{i \in A, j \in B} d_{ij}$. ■ **Average Linkage** (pairwise average of single and complete linkages) $d(A, B) = \frac{1}{n_A + n_B} \sum_{i \in A, j \in B} d_{ij} \rightarrow$ weighted average pairwise dissimilarity. ■ **Centroid Linkage**, $d(A, B) = \text{dist}(\bar{\mu}_A, \bar{\mu}_B)$

■ **Ward (Centroid),** based upon Euclidean distances \rightarrow weighted centroid linkage. $d(A, B) = \sqrt{\frac{n_A n_B}{n_A + n_B} \|\bar{\mu}_A - \bar{\mu}_B\|_2^2}$.

(21) ■ **Properties of Linkage:** ■ **Single Linkage** (1) Gives chaining, i.e., just add on single observation to a big cluster iteratively. (2) Chaining is flexible and good at handling different cluster shapes if the points are close together in their clustered shapes. (3) Good for outlier detection! If an observation has high height in the dendrogram, it means it was joined very late to the cluster, so it is very dissimilar or far from other observations. ■ **Complete Linkage:** (1) Optimal for spherical clusters/balls of uncorrelated features and balanced clusters. (2) Similar in properties to Kmeans and GMM however it is most robust to outliers! Whereas Kmeans is sensitive to outliers. ■ **Average Linkage:** In between single and complete. (1) Good for data that is close to spherical and balanced clusters. (2) Has statistical consistency - on average this gives the best clustering. ■ **Ward Linkage:** (1) Gives balanced and well-separated clusters. (2) ISSUE - can lead to inversion - flat parts of the dendrogram - bad property but tend to happen lower down in the tree so maybe it doesn't matter as much. ■ **NOTES:** When we have very correlated features all methods perform kind of bad - apply PCA first. ■ **NOTE:** THAT IN ORDER TO GET THE CLUSTER FOR SOME K WE "CUT" THE DENDROGRAM!!! DENDROGRAM - GOOD VIZ DATA SUMMARY.

■ **Bi-Clustering:** Cluster both the observations and the features. Can use a cluster heatmap - hierarchical clustering on features and observations separately and reorder the rows and columns according to the dendrogram order and plot the dendrogram on the x and y axis of the heatmap and fill in the heatmap according to some color gradient. We will see some good block like patterns if the data has groups of observations and features. Strong interpretations for both observations and features! ■ **Validation:** Clustering \rightarrow choosing K the number of clusters in the data. Goal: Minimize the within cluster distance/dissimilarity. Plot the within cluster distance across all values $K = 1, \dots, n$. For $K = 1$ the within cluster distance is high and decreases to $K = n$ where each observation is its own cluster and the within cluster distance is 0. This is like a training error curve and ends up not being useful for choosing the value of K . No simple way to get test error. So, this is a problem, and we need to try something else! ■ **Silhouette Statistic:** A measure of the goodness of cluster for each observation i . HIGHER IS BETTER \rightarrow most observations well clustered. Let $a_i =$ average within cluster dist for obs i and $b_i =$ average between cluster dist for obs i . The silhouette stat is $S_i = \frac{b_i - a_i}{\max(a_i, b_i)} \in [-1, 1]$. If $S_i \rightarrow 1$, then the average between cluster dist is large and the average within cluster dist is small hence the observation is well clustered. If $S_i \rightarrow -1$, then the observation is poorly clustered. Can choose $R = \text{argmax}_k \sum_{i=1}^n S_i$.

■ **ISSUES of Silhouette Stat:** (1) This yields strange behavior for $K = 1, 2$ and $K = n$, so useful if we try to determine if there are 3 to 10 clusters but can break if we are trying to decide 1 or 2 clusters. (2) Perform well for spherical and balanced clusters (okay for Kmeans, GMM, Complete, Ward's). (3) NEVER use this for single linkage, spectral clustering, or nonlinear dimensionality reduction technique + a method like Kmeans or GMM. Silhouette score assumes compact, convex clusters and uses Euclidean distances in the original space - this causes misleadingly low scores for methods like single linkage (which forms chained clusters), spectral clustering (which relies on graph-based similarity), or clustering after nonlinear dimensionality reduction like UMAP (which distorts distances and loses true structure). (4) Can help identify outliers too! ■ **Stability** (Consensus clustering): Idea: retain stable + reliable interpretation. Approach: repeated subsample of data and apply clustering with fixed K . Record the co-cluster membership \rightarrow Consensus Cluster Matrix $C_{n \times n}$ where entry $C_{ij} = \frac{\text{times } i \text{ and } j \text{ sampled together}}{\text{times } i \text{ sampled together}}$. If v_i was strong diagonal of I and off diagonal $b_{ij} =$ stable clustering across samples. Can compute variance of cluster assignment for observation i without C_{ij} $\rightarrow \frac{1}{n} \sum_{j=1}^n C_{ij} (1 - C_{ij})$. Then we can take $K = \text{argmin}_k \sum_{i=1}^n (1 - C_{ij})$ - choose K with smallest variance.

■ **Note:** we may have some clusters that are very stable on the diagonal as the within cluster variance is very low but the other clusters for K are less stable. This may indicate that if K is smaller, there is a cluster within a cluster (family/embedded cluster)!!!! Note that this stability method is done using bootstrap sampling to subsample! ■ **Adaptive Basis Function Models** (Supervised Learning): Regression $y = f(x) + \epsilon$. Now $f(x) = \beta_0 + \sum_{i=1}^n \beta_i \phi_i(x, \phi_0)$ where ϕ_0 is a basis function and ϕ_i are other basis functions and ϕ_0 are any parameters of the basis function. Here we are trying to learn weights β_i and basis function ϕ_i from the data. Such classes of supervised models are: (1) Decision Trees where ϕ_i are indicator functions (non-linear regression function). (2) Neural Networks where ϕ_i are other learned models (learn a model within a model).

(22) ■ **Decision Trees:** $\phi_n = \{x \in \mathbb{R}^d\}$, where \mathbb{R}^d is some rectangular region of our domain - fit + viz model as a binary rooted tree. ■ **Similar approach to hierarchical clustering** - locally greedy approximation (optimal decision at each step without regard to previous step) - start at top of tree and split data based on x . The optimal tree is NP-hard; therefore, we want a fast approximation using a GREEDY APPROACH \rightarrow Recursive Binary Partitioning (Divisive/Top-Down). Top-Down is faster computationally. ■ **Theory:** For each split: $\min_{f, f'} \left[\min_{\pi} \sum_{i \in \pi} L(y_i, f(x_i)) + \sum_{i \in \pi^c} L(y_i, f'(x_i)) \right] + \lambda \sum_{i \in \pi} L(y_i, f(x_i)) + \lambda \sum_{i \in \pi^c} L(y_i, f'(x_i))$. Want to minimize some type of loss where the observation is in this region vs it is not in this region. ■ **Example:** Loss \rightarrow Regression: $MSE(Y, \hat{y}) = \|y - \hat{y}\|_2^2$, prediction for region $\hat{\mu}_\pi = \frac{1}{n_\pi} \sum_{i \in \pi} y_i$ - sample mean for m^π region. ■ **Brute Force:** Try all π features and all possible split points. We can make this fast via a few tricks to reduce computation - faster than other ML algorithms! ■ **Classification:** $f \in \{1, \dots, K\}$. The loss $L(y, \pi, x) \rightarrow$ multi-class $= \hat{y}_k = \frac{\sum_{i \in \pi} \mathbb{1}(y_i = k)}{n_\pi}$ which is simply the proportion of observations in each class. Can use ■ **Misclassification Loss:** $\frac{1}{n_\pi} \sum_{i \in \pi} \mathbb{1}(y_i \neq \hat{y}_k)$ is fine to use but maybe not the best here as it does not lead to pure nodes! This is clearly not ideal.

■ **Cross Entropy Loss:** $\sum_{i \in \pi} \sum_{k=1}^K y_i \log(\hat{y}_{ik})$ (binary case $= \hat{y}_k \log \hat{y}_k + (1 - \hat{y}_k) \log(1 - \hat{y}_k)$ form for logistic regression). The issue is that it is slower computationally! ■ **Split Loss** (in between bins and CE): $\sum_{i \in \pi} \|y_i - \hat{y}_k\|_2$ is faster to update and leads to more pure nodes then misclassification loss. ■ **Size:** The largest possible tree is 2^n - pure node tree (training error is 0 - overfit). Can grow large tree and then "prune" (cross-complexity prune = regularization tree size). ■ **Properties:** (1) Trees can handle mixed types of features (continuous + binary/categorical + ordered). Don't need tons of preprocessing like one-hot encoding. (2) Missing values, can use as categorical feature, can use surrogate splits (mimics split), etc. (3) Interpretation of features using ■ **MDI (Mean Decrease in Impurity)** - difference in loss before the split and after the split - all splits. (4) Feature can be split on multiple times in a Decision Tree. (5) Bad for linear decision boundaries - struggle with lines! ■ **Results:** Decision Trees are not good predictors because they have high variance! The bias of tree depends on the size of the tree - if the tree is very deep, the bias will decrease and if the tree is deep, it will have high bias. If the decision boundary is linear the tree will incur a lot of bias. If the tree is large, the variance is very large! For moderate trees the variance is still high \rightarrow variance is always high. Weak learners \rightarrow poor predictors but better than random chance. Use many trees and LLN to reduce variance.

■ **Bagging** (Bootstrap aggregation): $x^{boot} = \text{bootstrap sample } x$ (w replacement). $f^{boot}(x) = \frac{1}{B} \sum_{b=1}^B f_b(x^{boot})$ - averaging over many f learners fit to many bootstrap samples. The bias stays the same and the variance is reduced! This is an ensemble method.

(23) ■ **Out-of-Bag Error (OOB Error):** Error of all ensemble members that didn't contain observation i evaluated at observation i , e.g., in a random forest is the error estimate computed by evaluating each training sample on the subset of decision trees for which that sample was NOT included in the bootstrap sample. Like LOOCV error. Built in validation method! To tune a random forest using OOB error, train multiple models with different hyperparameters and choose the one that performs best on data pointing it didn't see during training - as measured by its out-of-bag error.

■ **Theorem:** Bagging offers no reduction in variance for linear regression + classification models where the true underlying model is linear. Gauss-Markov Theorem. ■ **Properties:** These does bagging work? Works under non-linear models and methods. Typically, only used when fitting is very fast (or else computation is very expensive). ■ **Law of Large-Dimensions (LLN):** Assumes id. Are bagged trees are correlated. Therefore, to reduce variance we need to decorrelate. This is the idea behind RF. ■ **Random Forest:** Idea: Build an ensemble of decorrelated "random trees" (differ from Decision Trees because we only consider a random subset of features in the split). ■ **Algorithm:** For $b = 1, \dots, B$ - bootstrap sample x^{boot} . Build a random tree and at each split, consider a random subsample of features to split upon. Use the same tree building algorithm. Then take the ensemble average of all random forest trees. Note: RF is like an Adaptive KNN Algorithm. Average many decision boundaries = smooth. ■ **Notes:** Each random tree is a weak learner (non-optimal) but the idea, like ridge regularization, is that the BIAS increases (because we split on random subset of features), but the variance is reduced dramatically b/c we have a diverse ensemble of trees. NOTE: points should usually in terminal leaf with points that they are close to! ■ **Claim:** The random forest is the BEST of the shelf predictor for tabular data - don't have to do too much tuning, it just tends to work! Very hard to overfit. ■ **Hyperparameters:** (1) Tree depth: Large trees are favorable to drive down the bias (default = n). Features $B = \#$ of random trees. Use the OOB error (when it flattens out) to decide what value B . The smaller the tree depth, the more trees we need to fit! (3) m = max. features $= \#$ of features in subsample to split upon, if it is very large ($m \rightarrow n$) then this is just the same as bagging and if m is very small the performance will significantly degrade! The default $m = \text{max. features} = \sqrt{p}$.

■ **Feature Importance:** Take the average MDS over all RF Trees. RF is great for feature importance! ■ **Prediction:** $f^{boot}(x) = \frac{1}{B} \sum_{b=1}^B f^{boot}(x^{boot})$ for new datapoint we feed it through each of the RF tree predictors and take the average. ■ **Random Forest Advantages:** Great predictors (hard to overfit), strong interpretations but lose the tree-based interpretations of predictions, works for all data types and responses, great for missing data - Miss Forests \rightarrow based on missing data imputation. ■ **Random Forest Disadvantages:** (1) $p \gg n$ and we have a lot of irrelevant features, need $m = \text{max. features}$ to be large. (2) Highly correlated features have okay prediction, but it ruins MDS interpretation of the features!

■ **Intro to Model Stacking:** Idea: Diverse ensembles are great predictors! Ensemble members can be different types. Try to find the optimal f but why not take the ensemble of many w_f and take a weighted average (because some maybe better than others and should be weighted higher): $\min_L L(y; \sum_{f=1}^M w_f f(x)) = w_f f(x)$. We need to learn the weights w_f w/o new data to avoid overfitting - use training and validation to learn f and a QUERY set to learn the model stacking weights! (In RF $w_i = 1$).

(24) ■ **Boosting:** Sequentially/adaptively learn ensembles \rightarrow fitting to residuals; stage wise or greedy modeling. Intuition is to slowly fit residuals to hopefully outperform. ■ **AdaBoost:** $\eta_t(x) \in \{-1, 1\}$, $t = 1, \dots, M$. Algorithm: start with $f_0 = 0$ and weights $w_i = \frac{1}{n}$. Repeat for $m = 1, \dots, M$: (a) fit a weighted classifier $w(y_i, f_{m-1}(x_i))$ (b) weighted misclassification error: $\text{err}_m = \sum_{i=1}^n w_i \mathbb{1}(y_i \neq f_m(x_i)) / \sum_{i=1}^n w_i$. (c) Update ensemble coefficients $\hat{\alpha}_m = \log\left(\frac{1 - \text{err}_m}{\text{err}_m}\right)$. (d) Update weights: $w_i = w_i \exp(\hat{\alpha}_m \mathbb{1}(y_i \neq f_m(x_i)))$ and renormalize s.t. $\sum_{i=1}^n w_i = 1$. End loop. Output $= \text{sign}(\sum_{m=1}^M \hat{\alpha}_m f_m(x_i))$. ■ **Intuition:** (c) yields log odds of weighted misclassification error. If f_m is a good classifier the misclassification error is low and the $\hat{\alpha}_m$ is high (give more weight to the ensemble). (d) Suppose the i^{th} observation is misclassified and $\hat{\alpha}_m$ is large, then w_i increases significantly, if it is classified correctly then it is just $e^0 = 1$ so we keep the previous weight, but it normalizes so it down weights all correctly classified observations! FITTING TO MISCLASSIFIED POINTS AT EACH STAGE which is equivalent to fitting the residuals. Use Decision Trees as a base learner and shallow trees because they are very fast and weak learner \rightarrow weak learners make mistakes which is great for sequential learning. Avoid overfitting with small learning rate and weak learners! ■ **Forward-Stage-Wise Additive Modeling (FSAM)** (Greedy; Sequential learning): Repeat for $m = 1, \dots, M$. (a) $\hat{\mu}_m = \hat{h}_m(x) = \text{argmin}_L L(y; \hat{f}_{m-1}(x) + \beta h_m(x))$. (b) $\hat{f}_m(x) = \hat{f}_{m-1}(x) + \hat{\beta}_m \hat{h}_m(x)$. End loop. ■ **Intuition:** At each stage, fit base learner to the residual from the previous stage $L(y; \hat{f}_{m-1}(x))$. Not optimal because its fast to learn weight and learner! So can just fit weights as $\hat{\beta} =$ hyperparameter to slow it down, $\epsilon = 0$ - slow! ■ **Claim:** AdaBoost is FSAM with an exponential loss function (upper bound of CE). ■ **Epsilon Boosting** (slow boosting): FSAM where learning rate or weight is fixed to a small constant ϵ (hyperparameter λ). $\hat{f}_m(x) = \text{argmin}_L L(y; \hat{f}_{m-1}(x) + \epsilon h_m$