

Bootstrap Confidence Regions for Learned Feature Embeddings

Kris Sankaran

Department of Statistics
University of Wisconsin - Madison

ksankaran@wisc.edu

February 2, 2022

Abstract

Algorithmic feature learners provide high-dimensional vector representations for non-matrix structured signals, like images, audio, text, and graphs. Low-dimensional projections derived from these representations can be used to explore variation across collections of these data. However, it is not clear how to assess the uncertainty associated with these projections. We adapt methods developed for bootstrapping principal components analysis to the setting where features are learned from non-matrix data. We empirically compare the derived confidence regions in simulations, varying factors that influence both feature learning and the bootstrap. Approaches are illustrated on spatial proteomic data. Code, data, and trained models are released as an R compendium.

Dimensionality reduction is often used to explore the high-dimensional representations made by feature learning algorithms (Erhan, Bengio, Courville, & Vincent, 2009; Nguyen, Yosinski, & Clune, 2019). For example, principal components analysis (PCA) applied to word embedding vectors can highlight documents with similar content. These projections are typically treated as fixed, and their uncertainty is rarely characterized. A separate line of work, predating the widespread use of algorithmic feature learners, investigated bootstrap strategies for visualizing the uncertainty associated with dimensionality reduction methods (Josse, Wager, & Husson, 2016; Chateau & Lebart, 1996; Elguero & Holmes-Junca, 1988). The purpose of this study is to adapt these bootstrap methods to the algorithmic feature learning setting, accounting for difficulties that arise when using the output of one model as input to another (Wang, McCormick, & Leek, 2020; Chernozhukov et al., 2017).

For example, in the word embedding application above, training the embedding model twice may result in two different (though hopefully similar) representations of a document corpus. This introduces randomness in the feature extraction process that was not present when studying data-induced uncertainty in deterministic methods. Furthermore, the coordinates of the two learned embeddings need not correspond to one another.

This is different from the usual setting where, though new rows may be sampled from the population, the columns of the data matrix are assumed fixed. Moreover, many feature learning algorithms are supervised, and using labels to develop representations that are predictive of a response. If the model is overfit, the associated dimensionality reduction may give misleading interpretations, overstating the differences between classes in the learned feature space (Gross, Taylor, & Tibshirani, 2015; Friedman, Hastie, Tibshirani, et al., 2001).

For these reasons, visualizing uncertainty in learned feature embeddings requires additional care. We make the following contributions,

- Computing confidence areas: We describe three bootstrap-based strategies for computing confidence areas on dimensionality-reduced learned features. The approaches differ in the assumptions they make and the resources they require, and we clarify the effects of these choices on downstream analysis.
- Simulation and data analysis test-bed: We design generative mechanisms with enough complexity to warrant the application of feature learning algorithms, but which are transparent enough to support analysis. We prepare a dataset for a spatial proteomics problem where characterizing projection uncertainty has practical implications.
- Feature learner and dataset effects: We train supervised, unsupervised, or random-feature based learners using varying amounts of data and across degrees of model complexity. This clarifies whether projections from a supervised Convolutional Neural Network (CNN) require different treatment than those from an unsupervised Variational Autoencoder (VAE), for example.

Though all bootstrap approaches give comparable results and are approximately valid in a simple low-rank model, more complex simulations suggest that there is no universally applicable approach. In more complex contexts, an approach that only trains one feature learner appears to underestimate projection uncertainty across all types of feature extractors.

Section 1 summarizes technical tools used below. We present methods for constructing bootstrap confidence regions in Section 2, and we study its properties through simulation in Section 3. We compare methods in Section 4 through an application to a spatial proteomics dataset. We summarize in Section 5. We have released all code, raw and intermediate data, and trained models associated with simulations and data analysis. Documentation about these artifacts and how to reproduce them is provided in Supplementary Section 6.

1 Background

1.1 Feature Learning

We consider three complementary feature learning algorithms. The first is the VAE, which learns an L -dimensional reduction of a dataset by optimizing a reconstruction objective. It posits a generative model $p(z)p_\xi(x|z)$ of the data; $p(z)$ is a prior on latent features and $p_\xi(x|z)$ is a likelihood parameterized by ξ . The algorithm finds a pair ξ, φ maximizing the lower bound,

$$\log p_\xi(x) \geq \mathbb{E}_{q_\varphi} [\log p_\xi(x|z)] - D_{KL}(q_\varphi(z|x)||p(z))$$

where $q_\varphi(z|x) = \mathcal{N}(\mu_\varphi(x), \sigma_\varphi^2(x))$ maps raw data examples to distributions in a latent space. This problem is nonconvex, and the solution is non-deterministic. There are many implementations of VAEs; our experiments follow (Van Den Oord, Vinyals, et al., 2017).

Second, we learn supervised features through a CNN. A CNN regressor optimizes an empirical estimate of $\mathbf{E}\|y - f_{W_{1:J}}(x)^T \beta\|_2^2$ over $W_{1:J}$ and β . Here, $f_{W_{1:J}}$ transforms the raw input into the “final layer” features, and is defined recursively according to

$$\begin{aligned} f_{W_{1:j}}^j(x) &= \sigma\left(W_j f_{W_{1:(j-1)}}^{j-1}(x)\right) \\ f^0(x) &= x \end{aligned}$$

where $\sigma(x) := x\mathbb{I}(x \geq 0)$ and matrices W are restricted to the set of convolutions. Like in the VAE, this is solved through first-order optimization methods. Our implementation is the CBR architecture from (Raghu, Gilmer, Yosinski, & Sohl-Dickstein, 2017).

Third, we use a random convolutional features (RCF) model (Rahimi & Recht, 2008). A random sample of L training examples $x_{i_1}, \dots, x_{i_L} \in \mathbb{R}^{w \times h \times c}$ is selected; the x_i 's are assumed to be c -channel images with dimension $w \times h$. For each sample, a random $s \times s$ patch, denoted $w_p \in \mathbb{R}^{s \times s \times c}$, is extracted. For any c -channel image x , the l^{th} feature z_l is found by convolving x with w_l and spatially averaging over activations. This model uses random training image patches as convolutional kernels, rather than learning them from scratch. The features z_1, \dots, z_L are analogous to the features $f_{W_{1:J}}(x)$ in the CNN.

To train an RCF, the training data are featurized into $\mathbf{Z} \in \mathbb{R}^{n \times L}$. Then, a ridge regression model is trained from \mathbf{Z} to the y , giving an estimate $\hat{\beta}$. For a new example x^* , the same image patches w_1, \dots, w_L are used to form z^* , and predictions are made with $z^{*T} \hat{\beta}$. This model does not require gradient based training, and it can serve as a fast baseline.

1.2 Procrustes Analysis

Given centered \mathbf{X} and \mathbf{Y} , the Procrustes problem finds a rotation \mathbf{R} solving,

$$\min_{\mathbf{R} \in \mathcal{O}(p,p)} \|\mathbf{X} - \mathbf{Y}\mathbf{R}\|_F^2,$$

where $\mathcal{O}(p, p)$ is the space of $p \times p$ orthonormal matrices. The solution can be shown to be $\hat{\mathbf{R}} = \mathbf{U}^T \mathbf{V}$ for \mathbf{U} and \mathbf{V} obtained by the SVD of $\mathbf{X}^T \mathbf{Y}$ (Friedman et al., 2001; Gower, 1975). For B matrices $\mathbf{X}_1, \dots, \mathbf{X}_B$, the generalized Procrustes problem finds B rotations $\mathbf{R}_1, \dots, \mathbf{R}_B$ and mean \mathbf{M} solving

$$\min_{\mathbf{R}_1, \dots, \mathbf{R}_B \in \mathcal{O}(p, p), \mathbf{M}} \sum_{b=1}^B \|\mathbf{X}_b \mathbf{R}_b - \mathbf{M}\|_F^2.$$

While there is no closed form solution, the optimization can be solved by cyclically updating each \mathbf{R}_b via standard Procrustes problems and then updating $\mathbf{M} = \frac{1}{B} \sum_{b=1}^B \mathbf{X}_b \mathbf{R}_b$ (Friedman et al., 2001).

1.3 PCA and the Bootstrap

Several approaches are available for bootstrapping PCA. The total bootstrap computes B principal planes by applying PCA to B resampled versions of the data (Chateau & Lebart, 1996). For each replication, rows are sampled with replacement, viewed as draws from a larger population. The associated principal axes may be reflected or swapped with one another, so the associated sample coordinates are not directly comparable, and coordinates must be aligned. This is often accomplished through a Procrustes or conjoint analysis (Elguero & Holmes-Junca, 1988). In either case, the cloud of B points associated with each sample in the resulting reference space is used to form a confidence region for it.

In contrast, fixed-effects PCA views the rows of the data matrix as the entire population of interest (Josse et al., 2016). The source of randomness in this case is measurement noise around a low-rank model, not sampling from a larger population of rows. By fitting a measurement noise model and resampling residuals, a parametric bootstrap provides confidence regions for the true latent coordinates in the low-rank model. We also note that Bayesian factor analysis approaches sampling from the posterior of latent coordinates (Ren, Bacallado, Favaro, Holmes, & Trippa, 2017; Ren et al., 2020). Like the fixed-effects PCA model, these approaches specify an explicit low-rank model with measurement noise. Like the total bootstrap, underlying factors may be swapped, and alignment is necessary.

2 Methods

This section describes ways to adapt the bootstrap approaches above to the feature learning context. Our raw data are n samples $x_i \in \mathcal{X}$, where \mathcal{X} is the raw data domain, e.g., images, text sentences, or audio signals. A corresponding set $y_i \in \mathbb{R}$ of responses may be available. The full data are $\mathcal{D} = (x_i, y_i)_{i=1}^n$. A *feature learner* is a parameterized mapping $T(\cdot; \theta) : \mathcal{X} \rightarrow \mathbb{R}^L$ taking data from \mathcal{X} and representing it in \mathbb{R}^L .

For example, in a text data application, we expect the learner to transform a set of raw word sequences into a vector of features reflecting the topic of the document. θ is estimated from data, typically through an optimization,

$$\hat{\theta} := \arg \min_{\theta \in \Theta} \mathcal{L}(\mathcal{D}, T(\cdot; \theta)) \tag{1}$$

for some loss \mathcal{L} . In an unsupervised feature learner, candidates $\theta \in \Theta$ are functions of x_1, \dots, x_n alone. For a supervised feature learner, the class includes functions of both x_1, \dots, x_n and y . To simplify notation, we will write $z_i = T(x_i; \hat{\theta}) \in \mathbb{R}^L$ to denote the learned features for observation i .

A challenge is that the learned features are not the same from one run to the next; the l^{th} learned feature from run 1 need not have any relationship with the l^{th} feature from run 2. This is a consequence of using stochastic optimization in 1. However, even if there is no direct correspondence across runs, they may all reflect the same underlying latent features. In particular, projections from dimensionality reductions from the two datasets may be similar, after applying an appropriate alignment.

Suppose that data have been split into a feature learning set, indexed by $I \subset \{1, \dots, n\}$ and an inference set, indexed by I^C . The fraction $\frac{1}{n} |I|$ used for feature learning is a hyperparameter whose influence is empirically studied below. The learning set $(x_i)_{i \in I}$ is resampled B times, leading to B different feature extractors, $T(\cdot; \hat{\theta}^b)$ which can then be applied to the full dataset, yielding learned features $\mathbf{Z}_b \in \mathbb{R}^{n \times L}$

2.1 Nonparametric Bootstrap

Like the total bootstrap, one approach to comparing embeddings across feature extractors is to perform a dimensionality reduction on each and then align the projections. For each b , compute a singular value decomposition,

$$\mathbf{Z}_{b, I^C} = \hat{\mathbf{U}}_b \hat{\Sigma}_b \hat{\mathbf{V}}_b^T$$

where the index I^C means that only features associated with the inference set are used. Define coordinates for sample i with respect to the top K right singular vectors using $l_i^b = \sum_{k=1}^K \hat{u}_{ik}^b \hat{\sigma}_k^b$. These can be stacked into $\mathbf{L}_b \in \mathbb{R}^{|I^C| \times K}$.

A Procrustes analysis applied to $\mathbf{L}_1, \dots, \mathbf{L}_B$ learns a series of rotation matrices $\mathbf{R}_1, \dots, \mathbf{R}_B$ aligning the projections. For each sample i , compute a mean and covariance matrix based on the B vectors $\mathbf{R}_b l_i^b$. These are used to create $1 - \alpha$ level confidence areas for each inference sample in the K -dimensional projection space. This approach plugs-in an estimate for a “true” low-dimensional \mathbf{L} , assuming that this representation is noisily observed and then subject to arbitrary rotations. Note that if the true latent representations are subject to more general transformations (e.g., translations) across runs of the extractor, this assumption may not be appropriate.

The advantage of this approach is that it does not require a parametric model for simulating new versions of \mathbf{Z}_b . The price to pay is that it is necessary to train B feature extraction models $T(\cdot, \hat{\theta}_b)$, which can be a computational burden, even if it is parallelizable. Further, confidence areas are not computed for samples in the feature learning set $(x_i)_{i \in I}$. However, if the uncertainty of sample-level projections is assumed to vary smoothly, then a heuristic is to consider the uncertainty of a sample in I as comparable to those of its nearby samples in I^C .

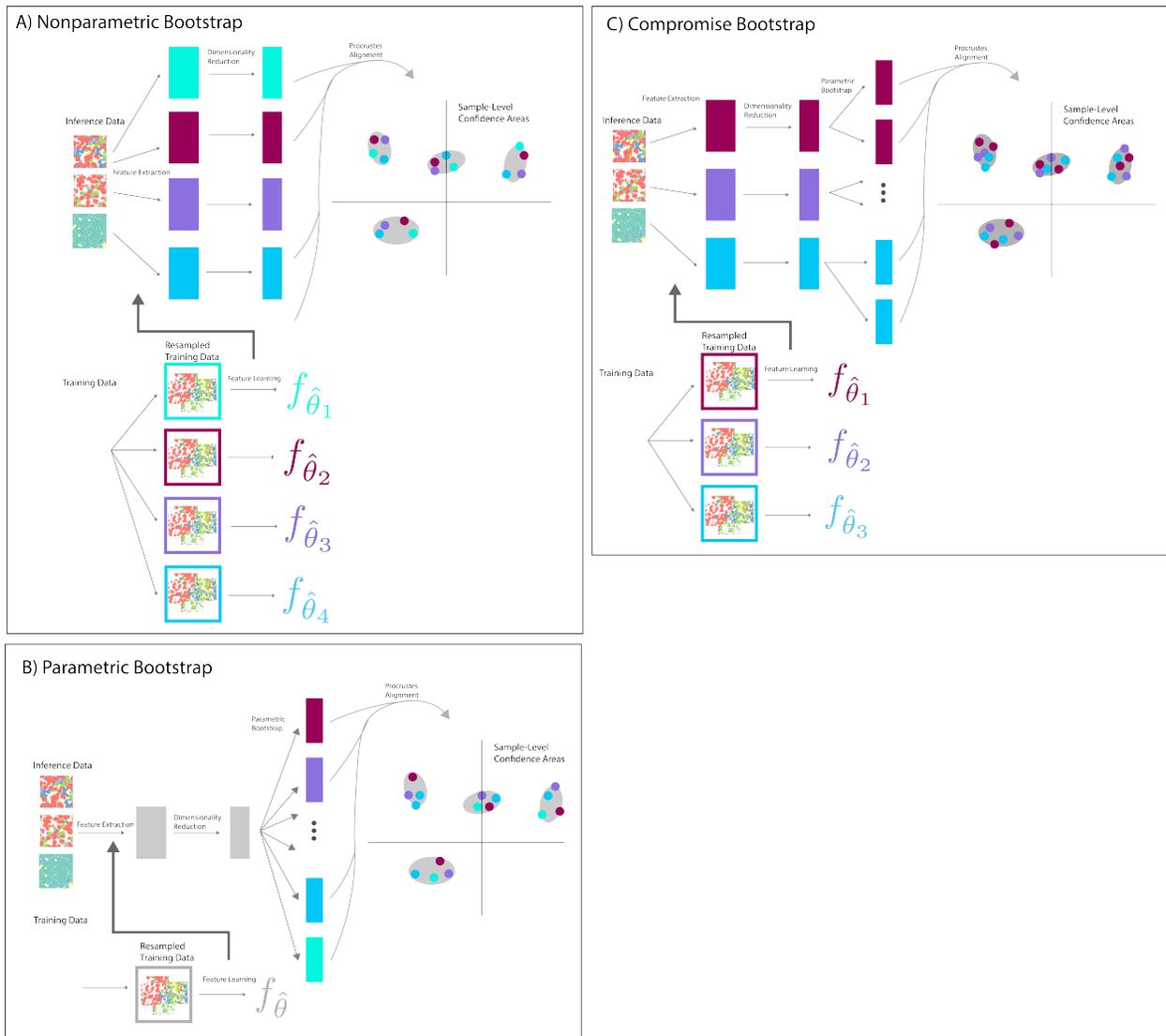


Figure 1: A summary of the proposed bootstrap procedures. All begin by splitting data into training and inference sets, used for feature learning and confidence region construction, respectively. The nonparametric bootstrap (a) trains B separate feature learners, each of which are used for feature extraction and dimensionality reduction before being aligned. The parametric bootstrap (b) trains a single feature learner and then simulates and aligns an ensemble of B latent coordinates for each sample. The compromise (c) trains a smaller set of feature learners but further resamples residuals (like in the parametric bootstrap) to increase the number of apparent bootstrap replicates.

2.2 Parametric Bootstrap

To avoid the computational complexity associated with training B feature extractors, we consider a parametric bootstrap, which simulates \mathbf{Z}_b by resampling residuals from an fitted low-rank model, analogous to the fixed-effects PCA approach (Josse et al., 2016). Suppose that variation across x_i is induced by latent features $l_i \in \mathbb{R}^K$. The feature learning process is modeled as,

$$\mathbf{Z} = \mathbf{L}\mathbf{V}^T + \mathbf{E} \quad E_{ij} \sim \mathcal{N}(0, \sigma_{\mathbf{E}}^2) \quad (2)$$

$$y = \mathbf{L}\beta + \epsilon \quad \epsilon_i \sim \mathcal{N}(0, \sigma_{\epsilon}^2) \quad (3)$$

where $\mathbf{L} \in \mathbb{R}^{n \times K}$ stacks the l_i and E_{ij} is the ij^{th} element of \mathbf{E} . Only \mathbf{Z} is available for predicting the response y .

To simulate \mathbf{Z}_b based on a single set of observed latent features \mathbf{Z} , we resample rows of \mathbf{Z} in the inference set I^C and compute the associated rank- K truncated SVD, $\hat{\mathbf{U}}\hat{\Sigma}\hat{\mathbf{V}}^T$. Then we draw,

$$\mathbf{Z}_b = \left(\hat{\mathbf{U}}\hat{\Sigma} + \mathbf{E}_b \right) \mathbf{\Pi}_b,$$

where $\mathbf{E}_b \in \mathbb{R}^{n \times K}$ is obtained by resampling entries of $\mathbf{Z} - \hat{\mathbf{Z}}$ and $\mathbf{\Pi}_b \in \mathbb{R}^{n \times K}$ is a random permutation matrix, reflecting the fact that coordinates of the feature extractor need not match from one run to the next. Alignment and confidence area construction then proceeds as in section 2.1.

2.3 Compromise

We adapt the mechanism above to simulate $\mathbf{Z}_1, \dots, \mathbf{Z}_B$ in the case where we have more than one trained feature extractor $T(\cdot, \hat{\theta}_s)$, for $s = 1, \dots, S$. Set $S < B$, so the feature learning phase is less costly than in section 2.1.

Begin by extracting \mathbf{Z}_s on resampled versions of the inference set, using the S extractors $T(\cdot, \hat{\theta}_s)$. Then compute their truncated, rank- K SVDs $\hat{\mathbf{U}}_s\hat{\Sigma}_s\mathbf{V}_s^T$. New feature sets are simulated from,

$$\mathbf{Z}_b = \left(\hat{\mathbf{U}}_{s(b)}\hat{\Sigma}_{s(b)} + \mathbf{E}_b \right) \mathbf{\Pi}_b,$$

where $s(b)$ is drawn uniformly from $1, \dots, S$ and \mathbf{E}_b resamples entries across all $\mathbf{Z}_s - \hat{\mathbf{Z}}_s$. Given the B resampled \mathbf{Z}_b , we generate confidence regions as before.

3 Simulations

We conduct two simulation studies. The first uses a low-rank model and permits calculation of coverage rates, but it is less representative of realistic feature learning settings. The second generates images using a spatial point process with variation reflecting a small set of latent parameters. The distributed feature learning associated with this setting prevents us from computing the coverage of confidence ellipsoids, but its complexity more accurately reflects practice.

3.1 Low-rank model

The first simulation generates samples $\mathbf{X} \in \mathbb{R}^{n \times D}$ using,

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T + \mathbf{E} \quad \mathbf{\Sigma} = \text{diag}(c\mathbf{1}_K) \quad \mathbf{U} \sim \text{Haar}(n, K) \quad (4)$$

$$y = \mathbf{U}\mathbf{\Sigma}\beta + \epsilon \quad \beta = \left(b\mathbf{1}_{\frac{K}{2}}, -b\mathbf{1}_{\frac{K}{2}}\right) \quad \mathbf{V} \sim \text{Haar}(D, K) \quad (5)$$

$$E_{ij} \sim \mathcal{N}(0, \sigma_{\mathbf{E}}^2) \quad (6)$$

$$\epsilon_i \sim \mathcal{N}(0, \sigma_y^2) \quad (7)$$

where $\text{Haar}(n, K)$ denotes a random orthonormal matrix with n rows and K columns. \mathbf{X} is a random rank- K matrix observed with Gaussian noise. y is a response depending on the latent coordinates of each row. The specific parameters we use are $N = 1000, D = 100, c = 100, b = 1, K = 2, \sigma_E^2 = 0.1$. Note that this is a model of the data \mathbf{X} , not the features \mathbf{Z} , as in equation 2.

As a feature extractor, we use a randomly perturbed and permuted SVD-based estimate of the latent coordinates,

$$\begin{aligned} \mathbf{X} &= \hat{\mathbf{U}}\hat{\mathbf{\Sigma}}\hat{\mathbf{V}}^T \\ \mathbf{Z} &:= \left(\hat{\mathbf{U}}_{\hat{K}}\hat{\mathbf{\Sigma}}_{\hat{K}} + \tilde{\mathbf{E}}\right)\mathbf{\Pi} \end{aligned} \quad (8)$$

where the subscript \hat{K} denotes that only the top \hat{K} left singular vectors and values are used and $\tilde{E}_{ij} \sim \mathcal{N}(0, 0.1^2)$. The permutation $\mathbf{\Pi}$ and noise $\tilde{\mathbf{E}}$ mimic the variation across retrained feature extractors. Given this source data and feature extractors, we apply all three bootstrap methods to this data, generating $B = 1000$ bootstrap replicates in each case. For the compromise approach, we train $S = 100$ separate feature extractors.

The resulting 95% confidence ellipses are shown in Figure 2. Qualitatively, the parametric and non-parametric bootstrap approaches provide similar output. A gradient across the colors of y reflects accurate estimation of the true latent factors. We have Procrustes aligned the true coordinates $\mathbf{U}\mathbf{\Sigma}$ (squares) with the B bootstrap replicates, and the fact that most squares are contained within ellipses suggests that the bootstrap accurately reflects uncertainty in the estimated projections. In fact, for the parametric and non-parametric approaches, the empirical coverage rates of these ellipses are 96.4% and 95.2%, respectively. On the other hand, the compromise approach appears to be overly conservative, with a coverage of 99.9%. This behavior arises in the remaining simulations and data analysis as well.

3.2 Spatial point process

In this simulation, we generate a collection of images using a point process where parameters vary from image to image. Intuitively, each image represents cells viewed through a microscope, and different latent parameters influence the cell ecosystem. A single response value y is associated with these latent parameters. Example images for varying y are given in Figure 3. We generate 10,000 of these $64 \times 64 \times 3$ -dimensional RGB images.

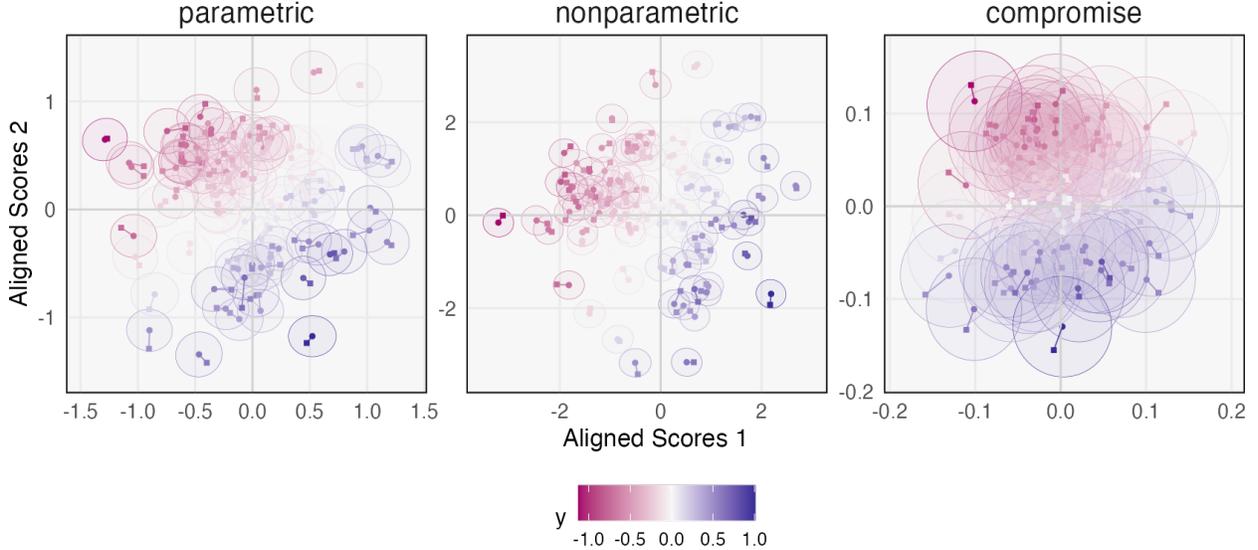


Figure 2: Projections from the low-rank data simulation. Each ellipse gives the confidence area for the latent coordinates of one sample. Squares are the positions of the true low-rank coordinates after Procrustes rotating them to align with the centers of the ellipses. The confidence areas for the nonparametric bootstrap are smaller than those for the parametric bootstrap. Those for the compromise method are conservative.

3.2.1 Generation

Locations of cells are governed by an intensity function drawn from a two-dimensional marked Log Cox Matern Process (LCMP) (Diggle, Moraga, Rowlingson, Taylor, et al., 2013). Recall that a Matern process is a Gaussian process with covariance function,

$$C_{\nu, \alpha}(\|x - y\|) = \sigma^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\sqrt{2\nu} \frac{\|x - y\|}{\alpha} \right)^\nu K_\nu \left(\sqrt{2\nu} \frac{\|x - y\|}{\alpha} \right), \quad (9)$$

where α acts like a bandwidth parameter and ν controls roughness.

Our LCMP has R classes (cell types). This can be constructed as follows. First, a nonnegative process $\Lambda(x)$ is simulated along the image grid, $\Lambda(x) \sim \exp(\mathcal{N}(0, \mathbf{C}_{\nu_\Lambda, \alpha_\Lambda}))$, where $\mathbf{C}_{\nu_\Lambda, \alpha_\Lambda}$ is the covariance matrix induced by equation 9. This is a baseline intensity that determines the location of cells, regardless of cell type. R further processes are then simulated, $B_r(x) \sim \exp(\beta_r + \mathcal{N}(0, \mathbf{C}_{\nu_B, \alpha_B}))$. These processes reflect relative frequencies of the R classes at any location x ; the intercept β_r makes a class either more or less frequent across all positions x . Given these intensity functions, we can simulate N cell locations by drawing from an inhomogeneous Poisson process with intensity $\Lambda(x)$. For a cell at location x , we assign it cell type r with probability $\frac{B_r^\tau(x)}{\sum_{r'=1}^R B_{r'}^\tau(x)}$. We have introduced a temperature τ controlling the degree of mixedness between cell types at a given location.

To complete the procedure for simulating images, we add two last sources of variation — the number of cells and cell size. The number of cells per image is drawn uniformly from 50 to 1000. The cells from class R are drawn with a random Gamma(5, λ_r) radius. A summary of all parameters used to generate each

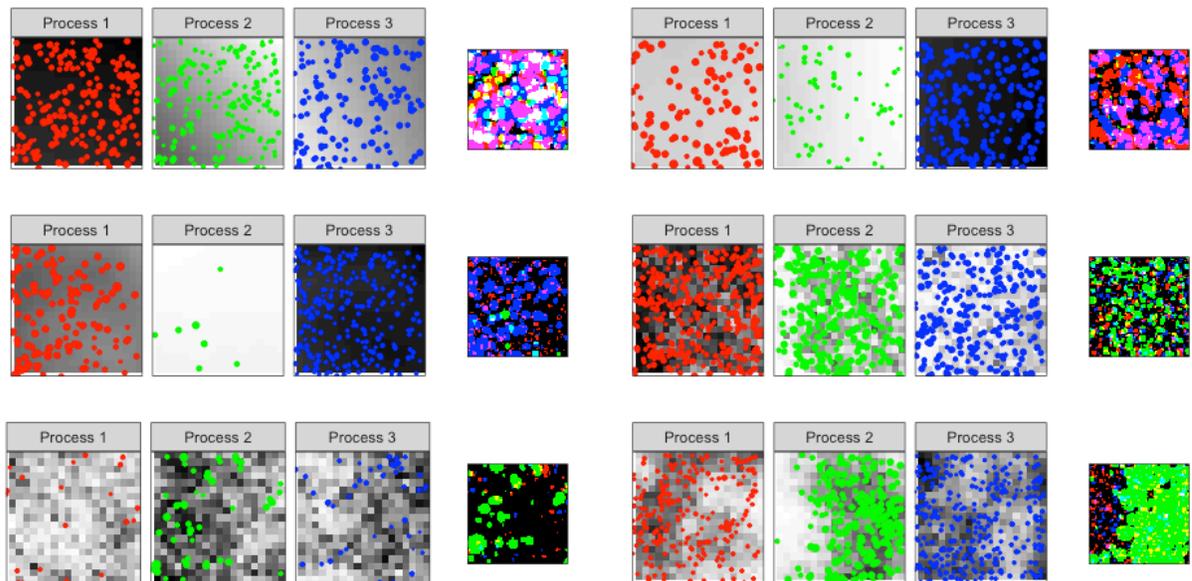


Figure 3: Example images, for low (top), average (middle), and high (bottom) values of y_i . For each sample, three relative intensity functions $B_r(x)$ are generated, shown as a greyscale heatmaps. Samples drawn from each process are overlaid as circles. The final images combine points across processes, removing the underlying intensity function, which is not available to the feature learner. Small y_i values are associated with smoother, less structured intensity functions.

image is given in Supplementary Table 1. Each parameter is drawn uniformly within its range, which has been chosen to provide sufficient variation in image appearance. These parameters are the “true” underlying features associated with the simulated images; they give the most concise description of the variation observed across the images. The response y is a hand-specified linear combination of these parameters, the reasoning behind the combination is discussed in the `generate.Rmd` script in the accompanying compendium, see Supplementary Section 6.

3.2.2 Experiments

We study the influence of the following parameters,

- Learning vs. inference split sizes. We vary the proportion of data used for learning and inference. We sample I so that $\frac{1}{n} |I| \in \{0.15, 0.5, 0.9\}$.
- Models trained. For feature extractors, we train CNN, VAE, and RCF models on the learning split I .
- Model complexity. We train VAEs whose hidden layer has dimensionality $L \in \{32, 64, 128\}$. Similarly, we vary the number of first-layer convolutional filters in the CNN model across $L \in \{32, 64, 128\}$. For the RCF, we use $L \in \{256, 512, 1024\}$ random features. This increase reflects the fact that more random features must be considered before a subset of predictive ones are identified.
- Inference strategy. We use the parametric, nonparametric, and compromise bootstrap strategies from section 2 to estimate confidence areas for the projections obtained by feature learners.

Figure 4 shows the activations of learned features across 2000 images for two perturbed versions of the training data when 90% of the data are used for inference and $L = 64$ (CNN, VAE) and 512 (RCF). The learned features correspond to,

- CNN: Activations from the final hidden layer of neurons, used directly as input for the regression.
- VAE: Spatially-pooled activations from the middle, encoding layer of the variational autoencoder.
- RCF: The spatially-pooled activations corresponding to each random convolutional feature.

Note that, across algorithms, there is no simple correspondence between learned and source features (i.e., parameters of the underlying simulation). Instead, there are clusters of learned features, corresponding to a pattern across multiple source features. We also find subsets of features across all models that are only weakly correlated with any source feature. This has been referred to as distributed representation learning (Hinton, 1984; Le & Mikolov, 2014).

Certain source features appear “easier” to represent than others, in the sense that more of the learned features are strongly correlated with them. Many features are correlated with N_i , the total number of cells in the image, and λ_{i1} , the size of the cells from Process 1. Depending on the model, the bandwidth α_{ir} ,

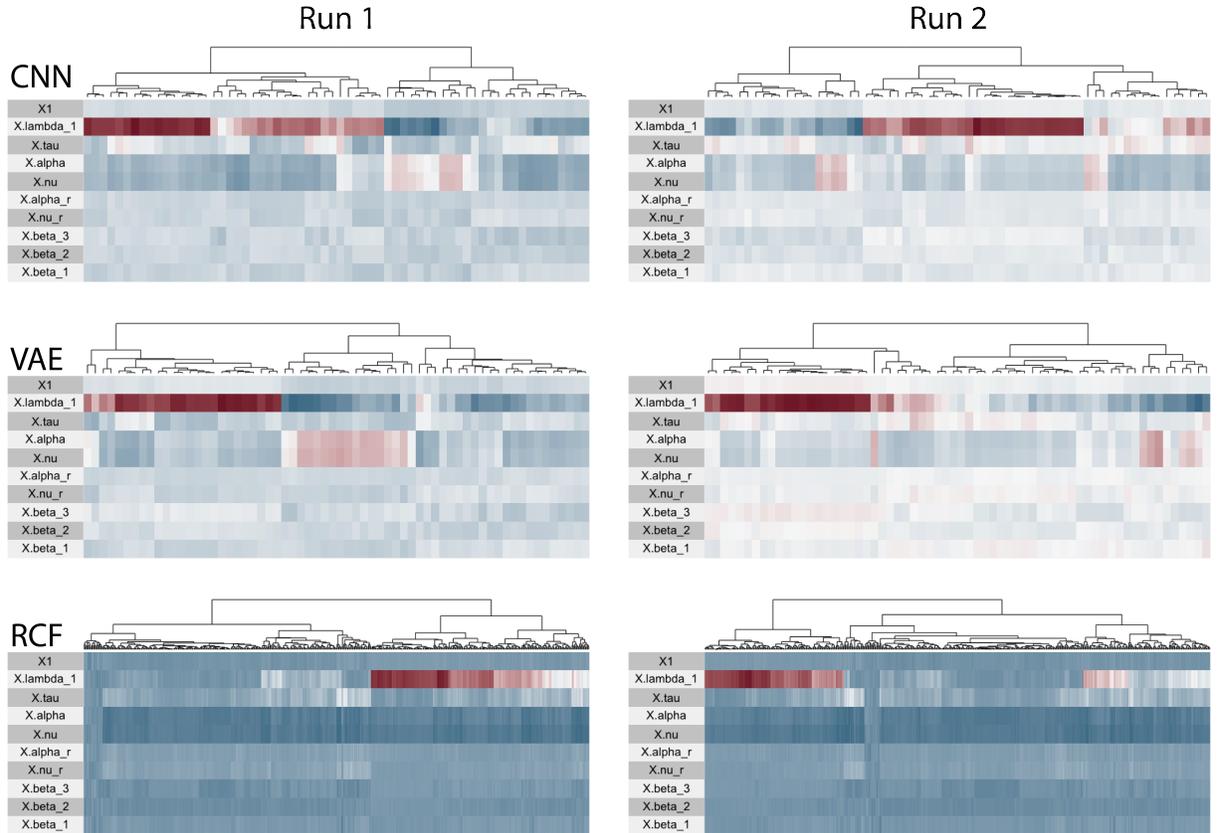
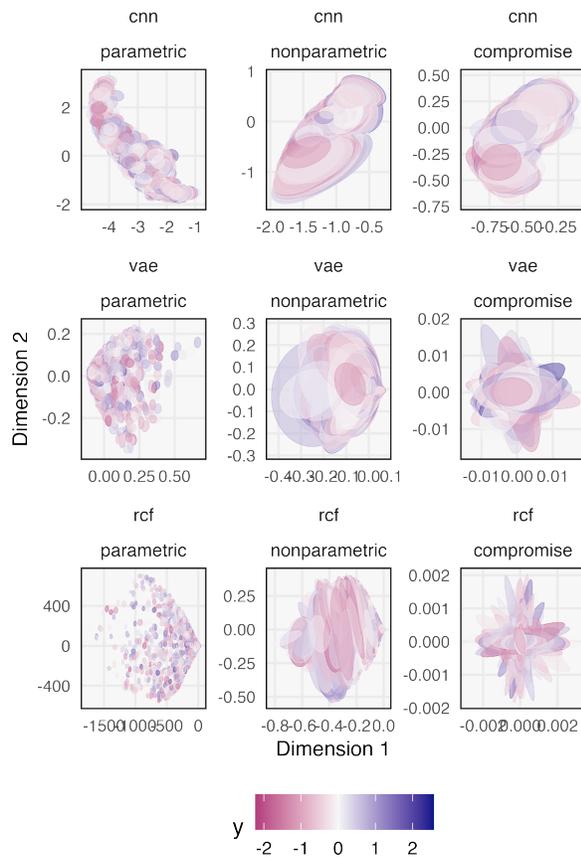


Figure 4: Each feature learning algorithm learns a distributed representation of the true underlying features in the simulation. Within each heatmap, rows correspond to the parameters in Supplementary Table 1. Columns are activations of learned features; they have been reordered using the package (Barter & Yu, 2018). The color of a cell gives the correlation between true and learned features. Blue and Burgundy encode positive and negative correlation, respectively.

roughness ν_{ir} , and prevalence β_{ik} parameters are either only weakly or not at all correlated with learned features. Even when features detect variation in α_{ir} and ν_{ir} , they cannot disambiguate between these two parameters. Finally, the CNN and VAE features tend to be more clustered, with strong correlation across several source features. In contrast, the RCF features show more gradual shifts in correlation strength. They also show relatively little variation in correlation strength across features other than λ_{i1} and N_i .

Example confidence areas across models and bootstrapping approaches are given in Figure 5a. In contrast to the Figure 2 of the low-rank simulation, the areas from the nonparametric bootstrap are larger than those from the parametric bootstrap. This disagreement suggests that the proposed mechanism of equations 4 and 8 is insufficient for characterizing differences in learned features that arise between runs of more complex feature extractors – multiple runs must be used for to account for randomness in algorithmic feature learning. As before, the compromise bootstrap has larger confidence areas than either parametric or nonparametric approaches on their own. In general, the RCF tends to have smaller confidence areas compared to the CNN

A. Varying learning and bootstrap methods



B. Varying splits and model complexity (CNN, parametric)

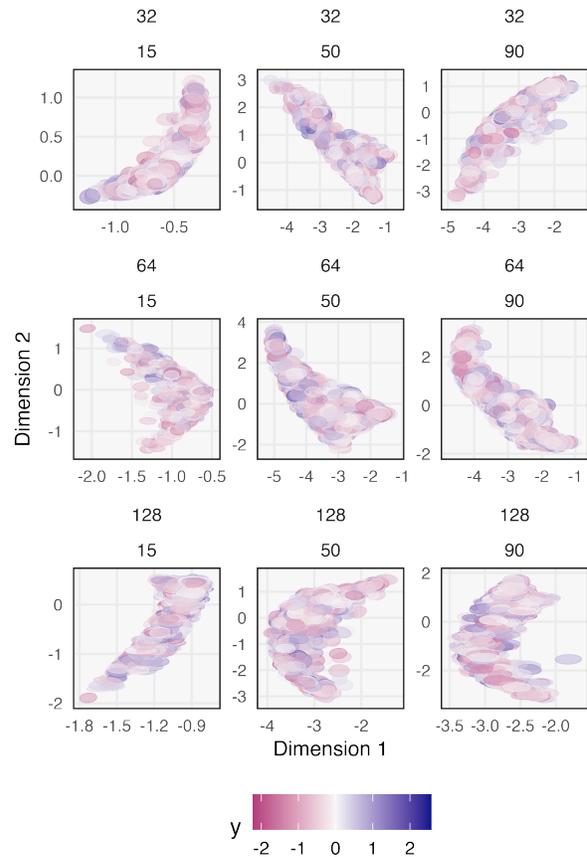


Figure 5: (a) The 95% confidence areas associated with projections from the spatial point process simulation. Each point correspond to one image. Only the setting with 90% of the data used for feature learning and the midsize models ($L = 64$ for the CNN and VAE, $L = 512$ for the RCF) are shown. (b) A view of confidence areas for the CNN across a range of learning split fractions (0.15, 0.5, 0.9) and model complexities ($L = 32, 64, 128$), all using the nonparametric approach.

and VAE.

Figure 5b shows confidence regions for a single model (CNN) and bootstrap procedure (parametric) across a range of model complexities and split proportions. For larger L , projections are further from the origin, suggesting larger activations on average. The fraction of data used for feature learning does not appear to affect the strength of the association with the response y or the size of the projection uncertainties. Corresponding figures for the other models are provided in Supplementary Section 7.

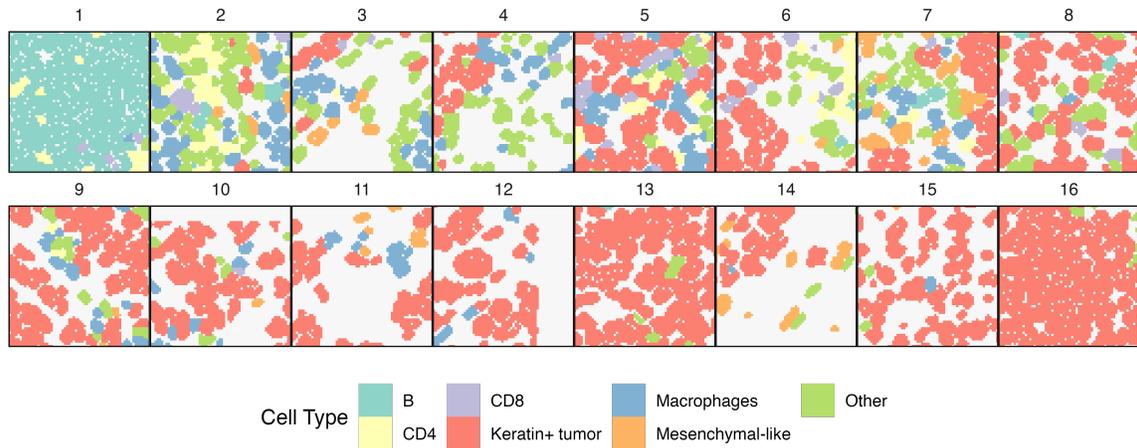


Figure 6: Example patches from the TNBC data. Panels are ordered by y_i , the (log) fraction of cells they contain that belong to the tumor. This provides signal for the supervised algorithms, whose goal is to correctly place patches from new patients along this gradient.

4 Data Analysis

We next analyze the spatial proteomics dataset reported in (Keren et al., 2018), which found a relationship between the spatial organization of Triple Negative Breast Cancer (TNBC) tissues and disease progression. In a classical proteomics study, the expression levels for a set of proteins is measured for a collection of cells, but the cell locations are unknown. In contrast, these data provide for each patient (1) an image delineating cell boundaries and (2) the protein expression levels associated with each cell in the images.

We only work with spatial cell delineations, not protein expression levels. This allows us to study feature learning within the images without having to worry about linking expression and image data, which is itself a complex integration problem. The data are 2048×2048 -dimensional images, one for each of 41 patients. Each pixel has a value 1 through 7 encoding which of 7 categories of tumor or immune cell types the pixel belongs. To ensure that the cell types are treated as categorical, we transform pixels to their one-hot encodings, resulting in a collection of $2048 \times 2048 \times 7$ binary matrices.

To setup a prediction problem, we split each image into $512 \times 512 \times 7$ patches. These patches are our x_i . Patches from 32 of the patients are reserved for feature learning. Four among these 32 are used as a development split, to tune parameters of the feature learning algorithms. As a response variable, we use $y_i = \log\left(\frac{\#\{\text{Tumor cells in } x_i\}}{\#\{\text{Immune cells in } x_i\}}\right)$. These y_i provide signal for the supervised feature learners. Example cell patches are shown in Figure 6. We fit the same models (CNN, VAE, and RCF) as discussed in section 3, varying model complexity over the same parameters as before.

As a baseline, we compare against a ridge regression with pixelwise composition features. We train a model with y as a response and the average number of pixels belonging to each of the cell-type categories as a 7-dimensional feature vector. This helps to determine whether the model has learned interesting features

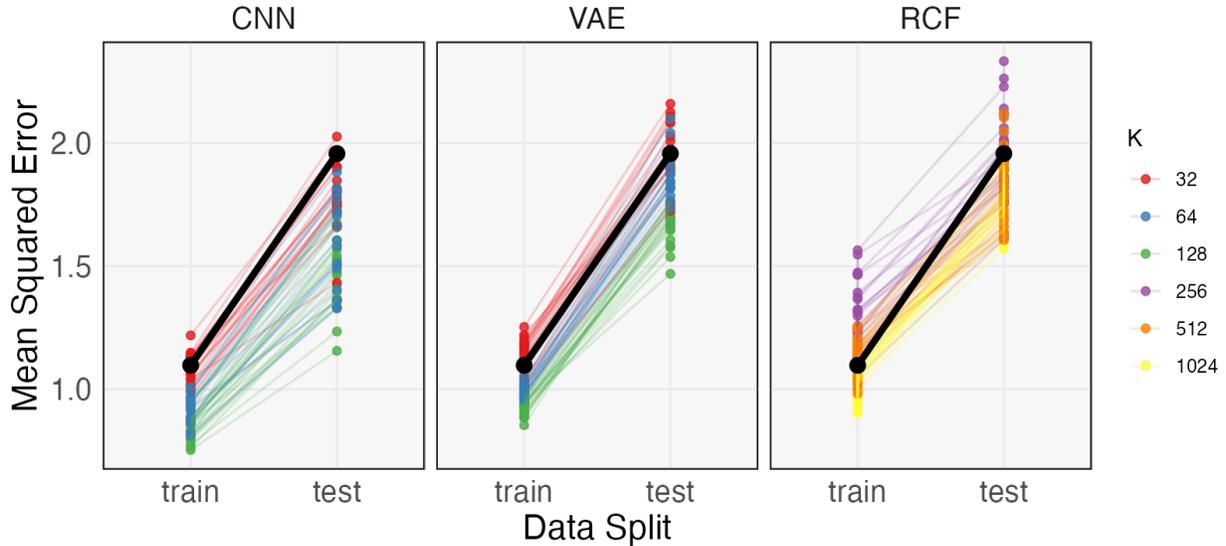


Figure 7: Relative performance of feature learning strategies on the TNBC data. Linked points come from the same bootstrap replicate. The black line gives the baseline ridge regression approach using the manually generated raw pixel counts for each of the cell types as predictors. Predictions from the development split are omitted; this split has few patches, and the estimated MSE’s have high variance.

for counting cells, like cell size and boundaries, rather than simply averaging across pixel values. Indeed, Figure 7 shows that, except in models with low capacity L , performance is improved when learning features algorithmically.

To characterize features and their uncertainty, we perform $B = 100$ iterations for each of the parametric, nonparametric, and compromise bootstrap strategies. In each case, the samples are generated from patches reserved for inference. Two-dimensional projections for fixed model complexities are given in Figure 8. As visible by the color gradients, all methods learn to differentiate between patches with small and large values of y_i , even the VAE, which is unsupervised. Comparing rows, the RCF appears to give the most stable representations, while the coordinates for the VAE have larger confidence areas in general. For all models, some projections appear to be more uncertain than others. Moreover, the certain axes directions tend to be more uncertain than others, reflected by the eccentricity of ellipses. For example, viewing estimates from the nonparametric approach, the VAE projections have the highest uncertainty for low values of Dimension 2. Analogously, high values of Dimension 2 have high uncertainty in the RCF.

For the CNN and RCF, the three bootstrap approaches give qualitatively similar conclusions about regions with higher and lower uncertainty, though the average sizes of the confidence areas differ. The size of confidence areas for the compromise approach in this case seem intermediate between those of the parametric and nonparametric approaches. For the VAE, the bootstrap approaches do not appear to agree. The compromise approach generally gives much larger confidence areas, potentially reflecting a failure of

the Procrustes alignment in this case. Though this figure only displays one L for each model, we find few differences in the projection uncertainty across models with different complexities, see Supplementary Figure 10 in the supplementary materials.

Figure 9 overlays example patches onto aligned coordinates. In the CNN, samples in the bottom right have a high fraction of immune cells, those on the top left are mostly tumor, and those at the top right have lower cell density. In light of the confidence regions in Figure 10, embeddings of immune cell rich images tend to show more variability across bootstraps. In the RCF, the first dimension similarly reflects the tumor vs. immune gradient. The lower uncertainty of projections along this axis suggests that this gradient is reliably captured across runs of the feature extractor. The lower right region of the VAE has high cell diversity, and the upper left has lower density. It seems that regions with higher cell diversity and larger proportions of immune cells also have more uncertain embeddings. From the top right to the bottom left, there is again a tumor to immune transition.

5 Discussion

We have adapted existing approaches for evaluating the uncertainty of projections in dimensionality reduction for use in the context of algorithmically learned features. We have conducted an empirical study using simulations of varying complexity and a spatial proteomics data analysis problem, applying a representative suite of feature learning algorithms. We found that in more complex settings, a parametric bootstrap based on a single set of learned features does not reflect the degree of uncertainty present when comparing features derived from independently trained models.

Our results raise several questions for further study. It is natural to ask to what extent similar behaviors are exhibited across other data domains, model types, or training regimes. For example, it would not be unusual to represent the cell data in our case study using a marked graph linking neighboring cells. Do the features learned by a graph autoencoder have similar stability properties? More generally, are there classes of feature learning algorithms that all share behavior from the stability perspective? In other domains, we may ask whether our methods can be adapted to text or audio data.

Though the proposed bootstraps provide similar conclusions in the low-rank simulation, they differ in the point process simulation and spatial proteomics data analysis. This suggests that mechanisms in equations 4 and 8 may not reflect the behavior of repeated feature learning in more complex situations. The assumption that features can be rotated to align runs may also be problematic, and more general transformations across feature learning runs are plausible. For the CNN and RCF models in the data analysis example, we find that confidence areas for the compromise approach are intermediate between those for the parametric and nonparametric bootstraps. However, in both simulations, it tends to be larger, and we have no explanation. Finally, though we have empirically found coverage rates for the parametric and nonparametric bootstraps to be acceptable in the low-rank simulation, we have not theoretically studied the properties of these procedures.

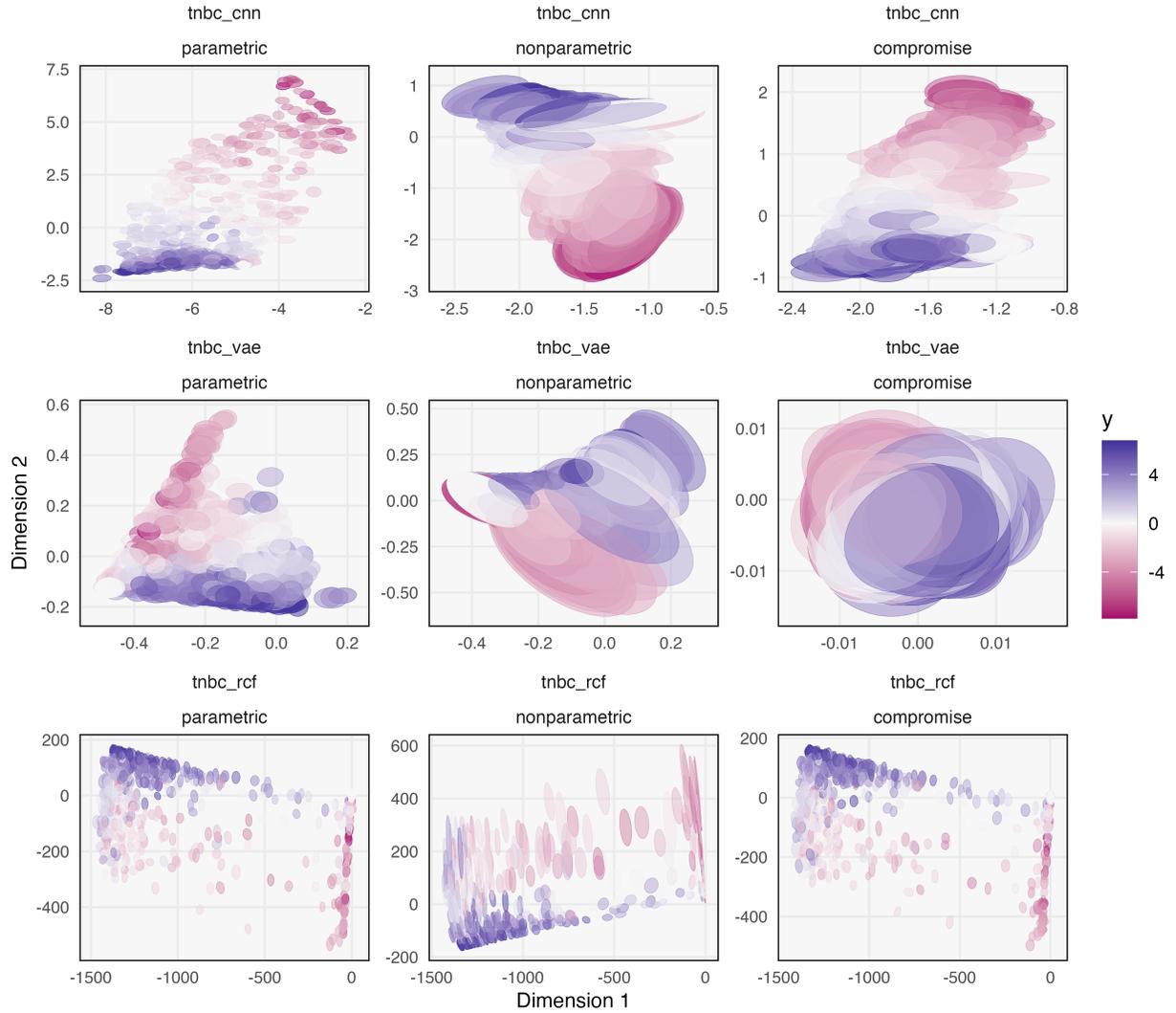


Figure 8: Confidence areas from the TNBC application. Points are shaded by $y_i = \log \left(\frac{\#\{\text{Tumor cells in } x_i\}}{\#\{\text{Immune cells in } x_i\}} \right)$, which provides the supervisory signal to the CNN and RCF during feature extraction. Models and bootstrap procedures are arranged along rows and columns, respectively. Only the models with intermediate complexity ($L = 64$ for the CNN and VAE, $L = 512$ for the RCF) are shown. Analogous figures for other L are given in the supplementary materials.

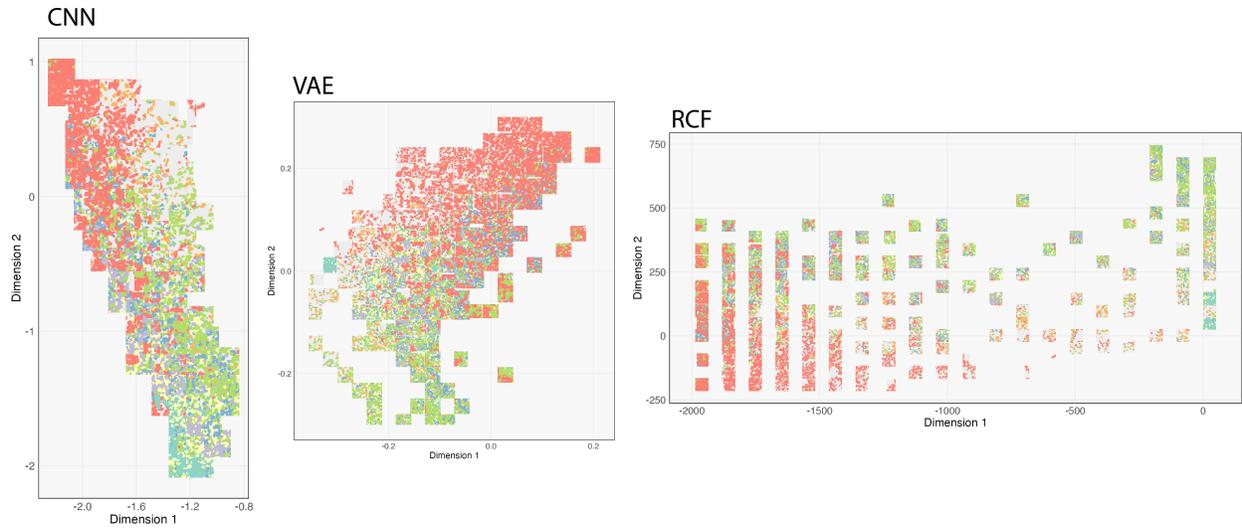


Figure 9: A version of the nonparametric bootstrap column from the dimensionality reductions in Figure 8, overlaying representative samples across the learned feature space. Cells are color coded as in Figure 6. Note that the overall shape of the region in which images are displayed mirrors the shape of the cloud of points in Figure 8.

This could clarify differences that arise in projection uncertainties between bootstrap and feature learning approaches.

Acknowledgments

The author thanks Susan Holmes, Karl Rohe, three reviewers, and the editor for feedback which improved the manuscript. Research was performed with assistance of the UW-Madison Center For High Throughput Computing (CHTC).

References

- Barter, R. L., & Yu, B. (2018). Superheat: An R package for creating beautiful and extendable heatmaps for visualizing complex data. *Journal of Computational and Graphical Statistics*, 27(4), 910–922.
- Chateau, F., & Lebart, L. (1996). Assessing sample variability in the visualization techniques related to principal component analysis: bootstrap and alternative simulation methods. In *Compstat* (pp. 205–210).
- Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., & Newey, W. (2017). Double/debiased/neyman machine learning of treatment effects. *American Economic Review*, 107(5), 261–65.

- Diggle, P. J., Moraga, P., Rowlingson, B., Taylor, B. M., et al. (2013). Spatial and spatio-temporal log-gaussian cox processes: extending the geostatistical paradigm. *Statistical Science*, 28(4), 542–563.
- Elguero, E., & Holmes-Junca, S. (1988). Confidence regions for projection pursuit density estimates. In *Compstat* (pp. 59–63).
- Erhan, D., Bengio, Y., Courville, A., & Vincent, P. (2009). Visualizing higher-layer features of a deep network. *University of Montreal*, 1341(3), 1.
- Friedman, J., Hastie, T., Tibshirani, R., et al. (2001). *The Elements of Statistical Learning* (Vol. 1) (No. 10). Springer series in statistics New York.
- Gower, J. C. (1975). Generalized procrustes analysis. *Psychometrika*, 40(1), 33–51.
- Gross, S. M., Taylor, J., & Tibshirani, R. (2015). A selective approach to internal inference. *arXiv preprint arXiv:1510.00486*.
- Hinton, G. E. (1984). Distributed representations.
- Josse, J., Wager, S., & Husson, F. (2016). Confidence areas for fixed-effects pca. *Journal of Computational and Graphical Statistics*, 25(1), 28–48.
- Keren, L., Bosse, M., Marquez, D., Angoshtari, R., Jain, S., Varma, S., ... others (2018). A structured tumor-immune microenvironment in triple negative breast cancer revealed by multiplexed ion beam imaging. *Cell*, 174(6), 1373–1387.
- Le, Q., & Mikolov, T. (2014). Distributed representations of sentences and documents. In *International conference on machine learning* (pp. 1188–1196).
- Nguyen, A., Yosinski, J., & Clune, J. (2019). Understanding neural networks via feature visualization: A survey. In *Explainable ai: interpreting, explaining and visualizing deep learning* (pp. 55–76). Springer.
- Raghu, M., Gilmer, J., Yosinski, J., & Sohl-Dickstein, J. (2017). Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability. In *Advances in neural information processing systems* (pp. 6076–6085).
- Rahimi, A., & Recht, B. (2008). Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning. *Advances in neural information processing systems*, 21, 1313–1320.
- Ren, B., Bacallado, S., Favaro, S., Holmes, S., & Trippa, L. (2017). Bayesian nonparametric ordination for the analysis of microbial communities. *Journal of the American Statistical Association*, 112(520), 1430–1442.
- Ren, B., Bacallado, S., Favaro, S., Vatanen, T., Huttenhower, C., & Trippa, L. (2020). Bayesian mixed effects models for zero-inflated compositions in microbiome data analysis. *The Annals of Applied Statistics*, 14(1), 494–517.
- Van Den Oord, A., Vinyals, O., et al. (2017). Neural discrete representation learning. In *Advances in neural information processing systems* (pp. 6306–6315).
- Wang, S., McCormick, T. H., & Leek, J. T. (2020). Methods for correcting inference based on outcomes predicted by machine learning. *Proceedings of the National Academy of Sciences*, 117(48), 30266–

6 Reproducibility

We have prepared a research compendium at <https://github.com/krisrs1128/LFBCR>. Instructions for downloading raw and preprocessed data are given in the `analysis/data/` subdirectory. The LFBCR compendium available in the repository above is also an R package and includes functions for that can be used to apply the bootstrap and alignment methods more generally.

Rmarkdown files needed to reproduce the figures in Section 3 are given in the `analysis/simulations` subfolder. The analogous files for the spatial proteomics application are given in `analysis/data_analysis`. The `learning` subfolders at both these locations can be used to retrain one instance of each of the feature learning algorithms. Further details are given in the `README.md` folders in the repository.

7 Supplementary Tables and Figures

| Feature | Description | Influence | Range |
|----------------------|--|------------------------------|---------------|
| N_i | The total number of cells. | 0.5 | [50, 1000] |
| $\nu_{i,\Lambda}$ | The roughness of the overall intensity process. | -0.5 | [0, 8] |
| $\alpha_{i,\Lambda}$ | The bandwidth of the overall intensity process. | -0.5 | [0, 8] |
| β_{ir} | The intercept controlling the frequency of class r . | 1 for $r = 1$, -1 otherwise | [-0.15, 0.15] |
| ν_{iB} | The roughness of the relative intensity processes. | -0.5 | [0, 3] |
| α_{iB} | The bandwidth of relative intensity processes. | -0.5 | [0, 3] |
| τ_i | The temperature used in cell type assignment. | 0.5 | [0, 3] |
| λ_{ir} | The shape parameter controlling the sizes of each cell type. | 1 for $r = 1$, 0 otherwise | [100, 500] |

Table 1: Our simulation mechanism is governed by the above parameters. Parameters N_i and λ_{ir} control the number and sizes of imaginary cells. $\nu_{i,\Lambda}$, $\alpha_{i,\Lambda}$, β_{ik} , ν_{iB} , α_{iB} , and τ_i control the overall and relative intensities of the marked LCMP from which the cell positions are drawn. Example draws are displayed in Figure 3.

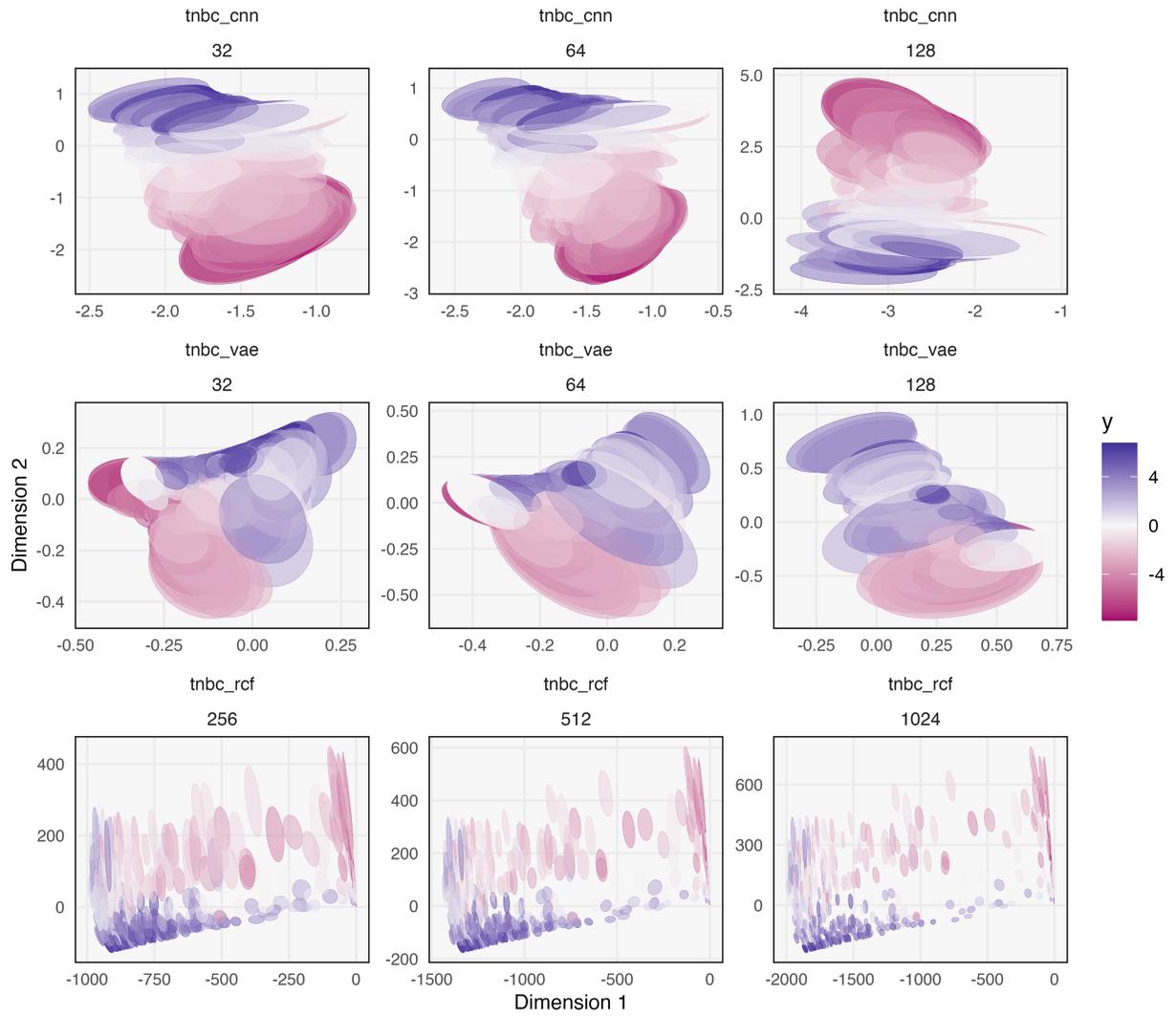


Figure 10: A comparison of the influence of models (rows) and model complexities (columns) on confidence regions constructed using the nonparametric bootstrap in the data analysis application. Each individual panel is read similarly to Figure 8.

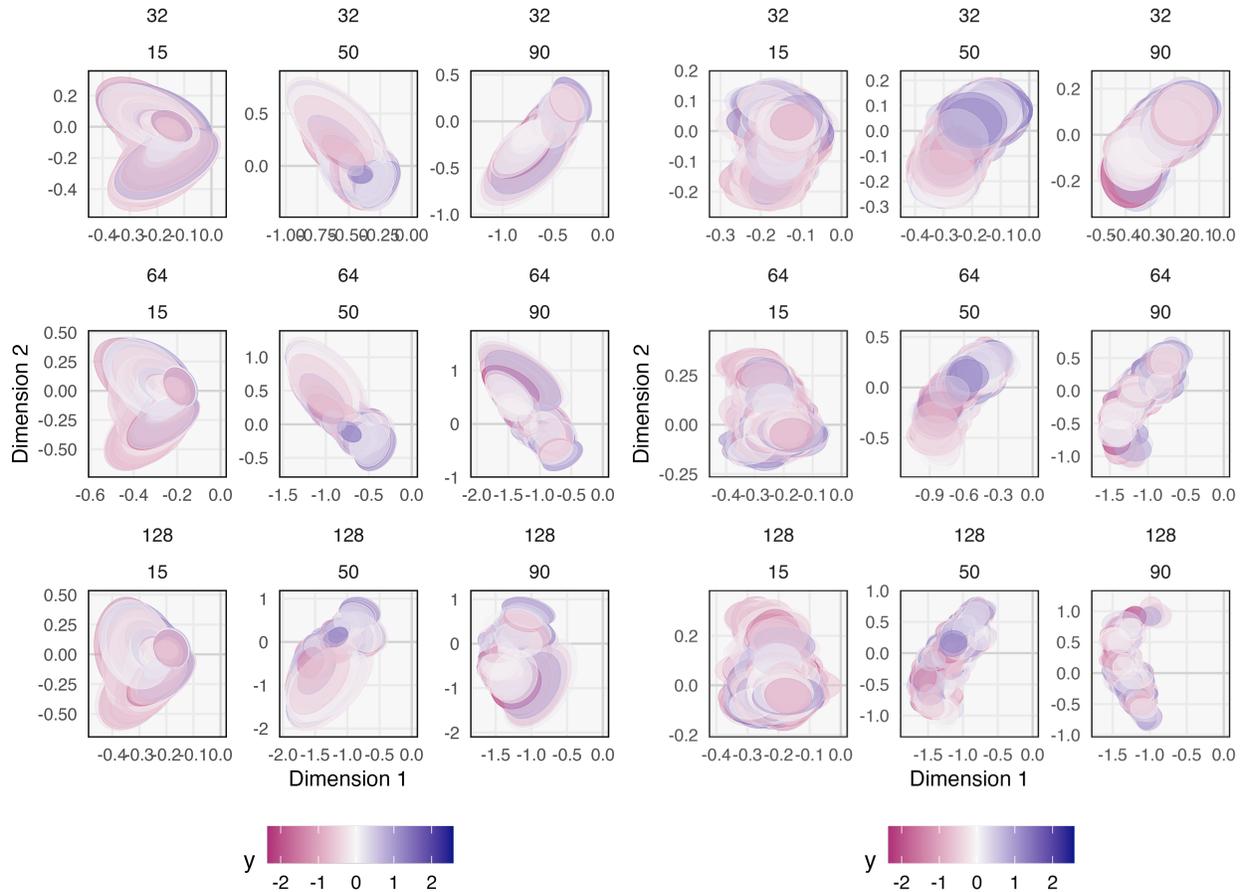


Figure 11: The analog of Figure 5, still shown on CNN projections, but for the nonparametric (left) and compromise (right) bootstrap approaches.

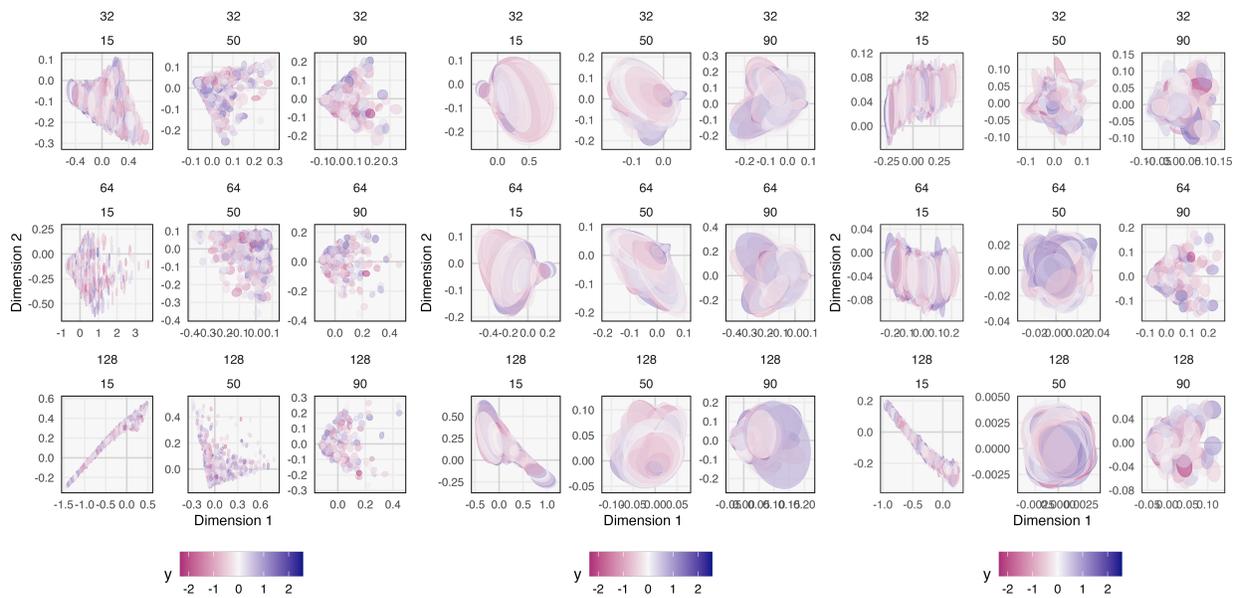
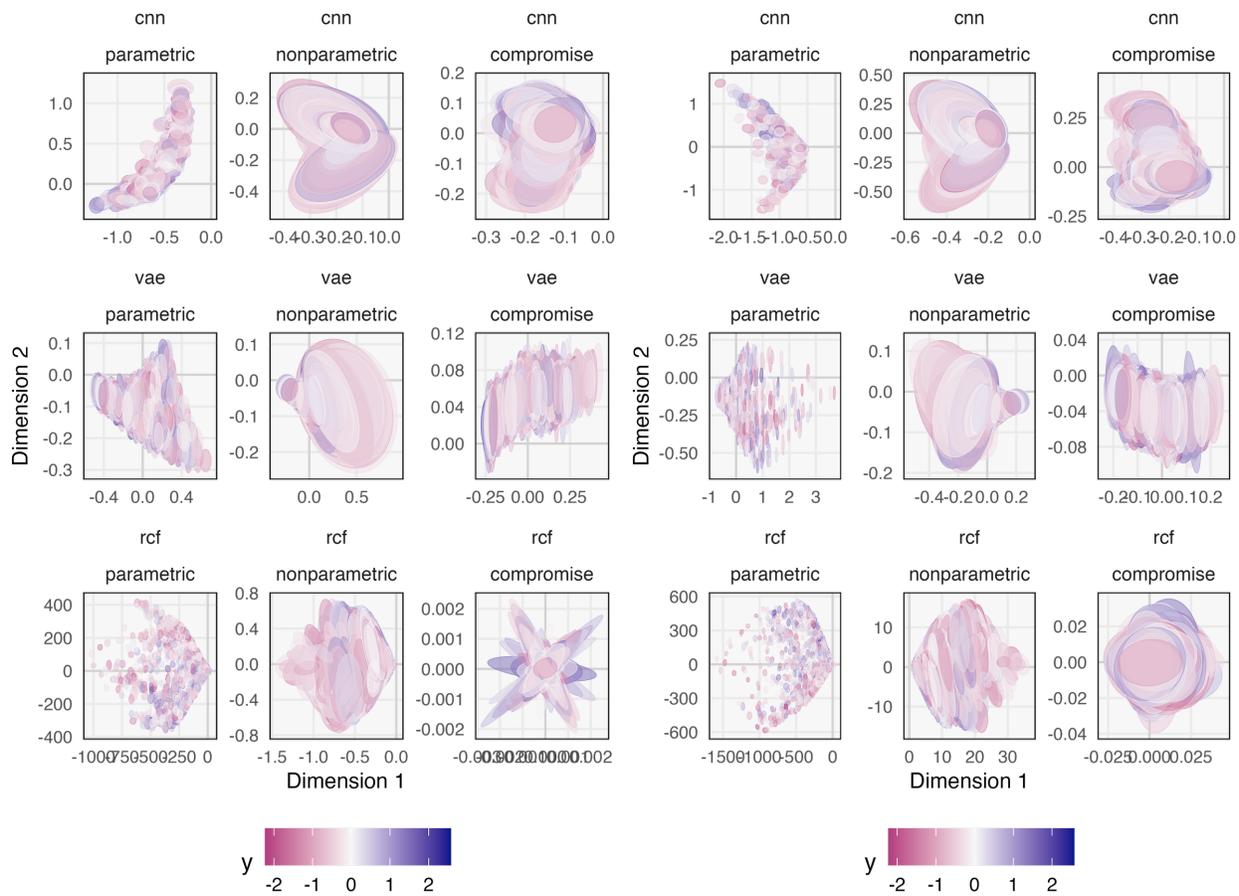


Figure 12: The analog of Figure 5, but for VAE models evaluated using the parametric (left), nonparametric (middle) and compromise (right) bootstrap approaches.



(b)

Figure 13: Confidence areas corresponding to Figure 5a, but for models with lower and average complexity, still using 15% of data reserved for learning.

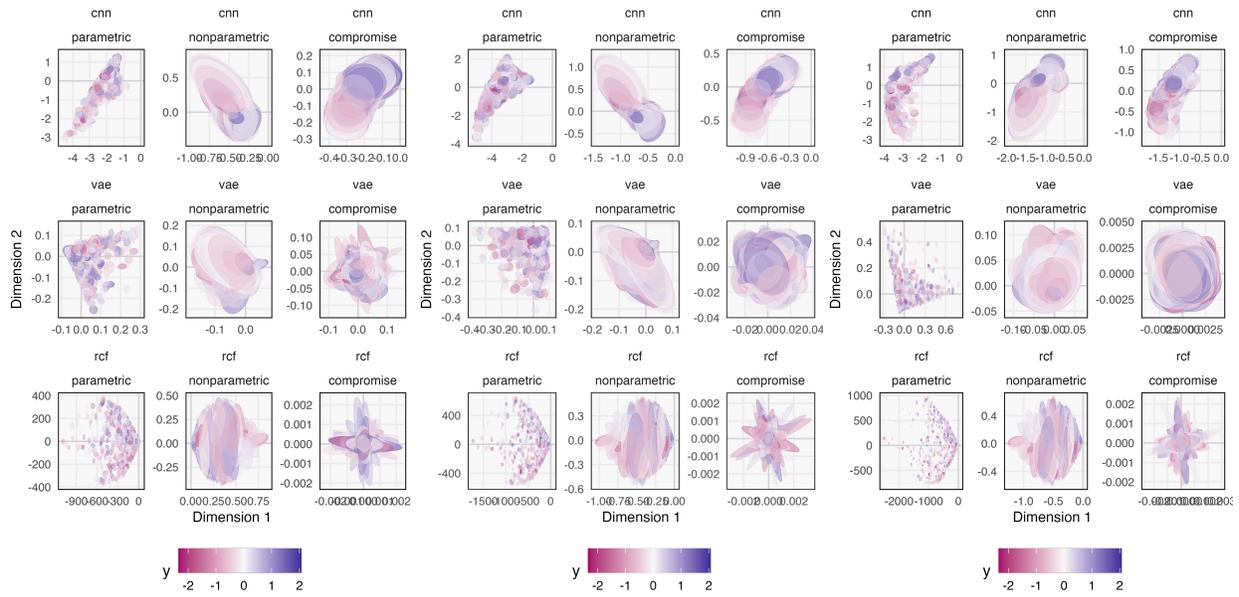


Figure 14: Confidence areas analogous to Figure 5a, but for models trained using 50% of the data. The source images can be viewed in full resolution at the compendium repository.

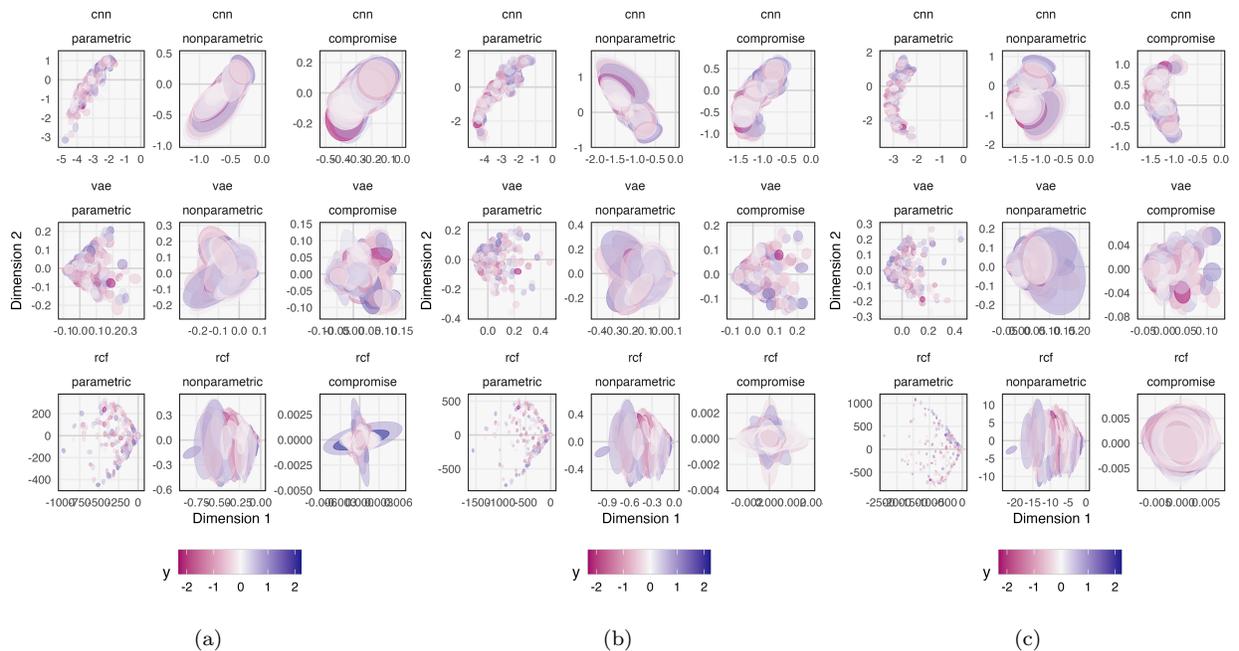


Figure 15: Confidence areas analogous to Figure 5a in the spatial point process simulation, but for models trained using 90% of the data. The source images can be viewed in full resolution at the compendium repository.

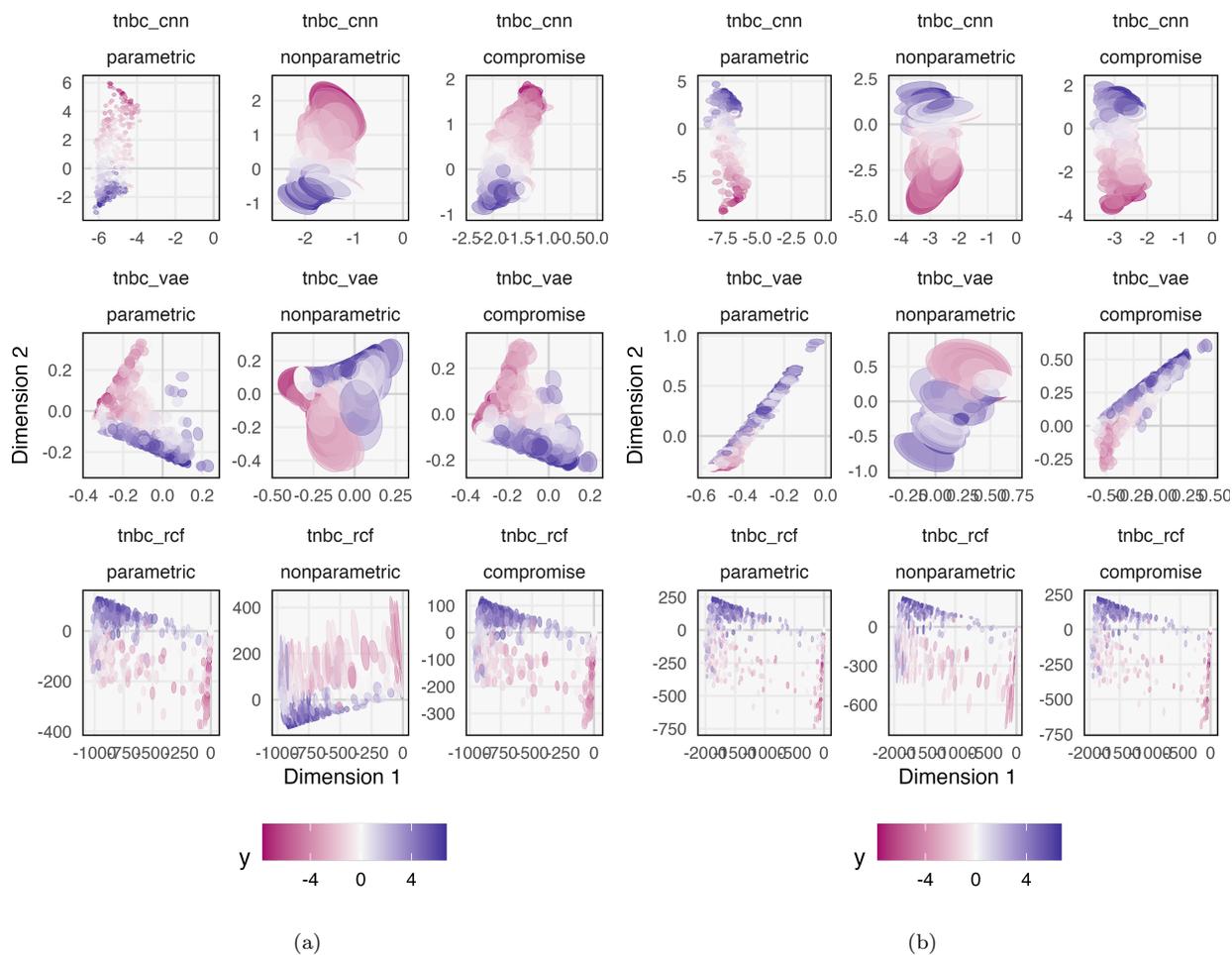


Figure 16: Confidence areas analogous to Figure 8 in the data analysis application, but for models with lower (left) and higher (right) complexity. The source images can be viewed in full resolution in the compendium repository.