

# Parzen-PNN Gaussian Mixture Estimator: Results Report

Fabrizio Benvenuti

February 3, 2026

## 1 Results

### 1.1 Parzen Window

In this section are reported the estimation results ValNLL and MSE obtained by varying the input parameters of the Parzen Window estimator.  
(i.e., window size  $h_1$  and number of sampled points per gaussian in the mixture).

### 1.1.1 Parzen Window Errors

Parzen Window Errors (Mixture 1) — normalized z

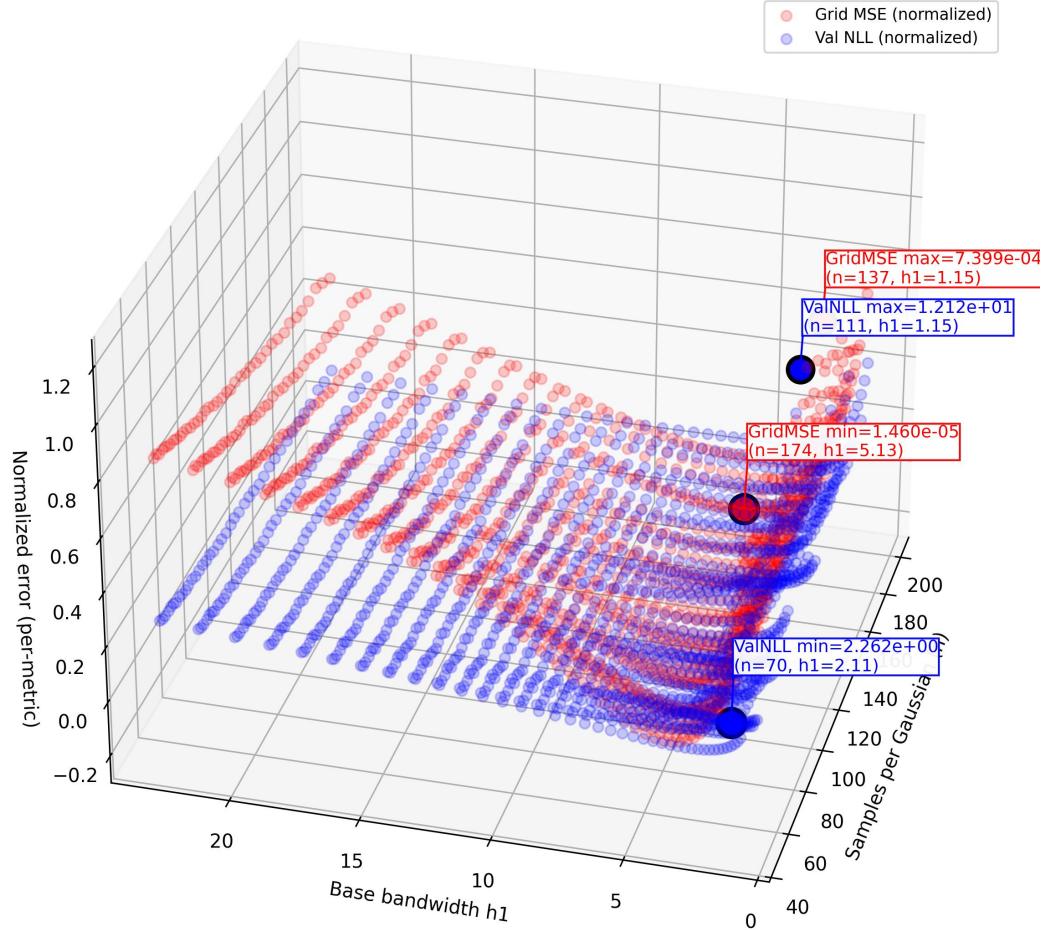


Figure 1: In red: MSE between the mixture 1 pdf and its estimate PW estimate  
 In blue: ValNLL between PW estimate and the sampled points;  
 while varying the window size  $h_1$  and the number of sampled points for each gaussian.

In this graph it's also noticeable how the MSE does not vary linearly with window size; approaching zero at the optimal value of  $h_1$  and increasing exponentially when undersmoothing occurs; transitioning to the oversmoothed region, the MSE still increases but less steeply.

Samples per gaussian seem to have an exponential impact on reducing MSE, even if its effect is less visible than the base bandwidth one.

NLL does not show the same steepness in the oversmoothing region as MSE does.

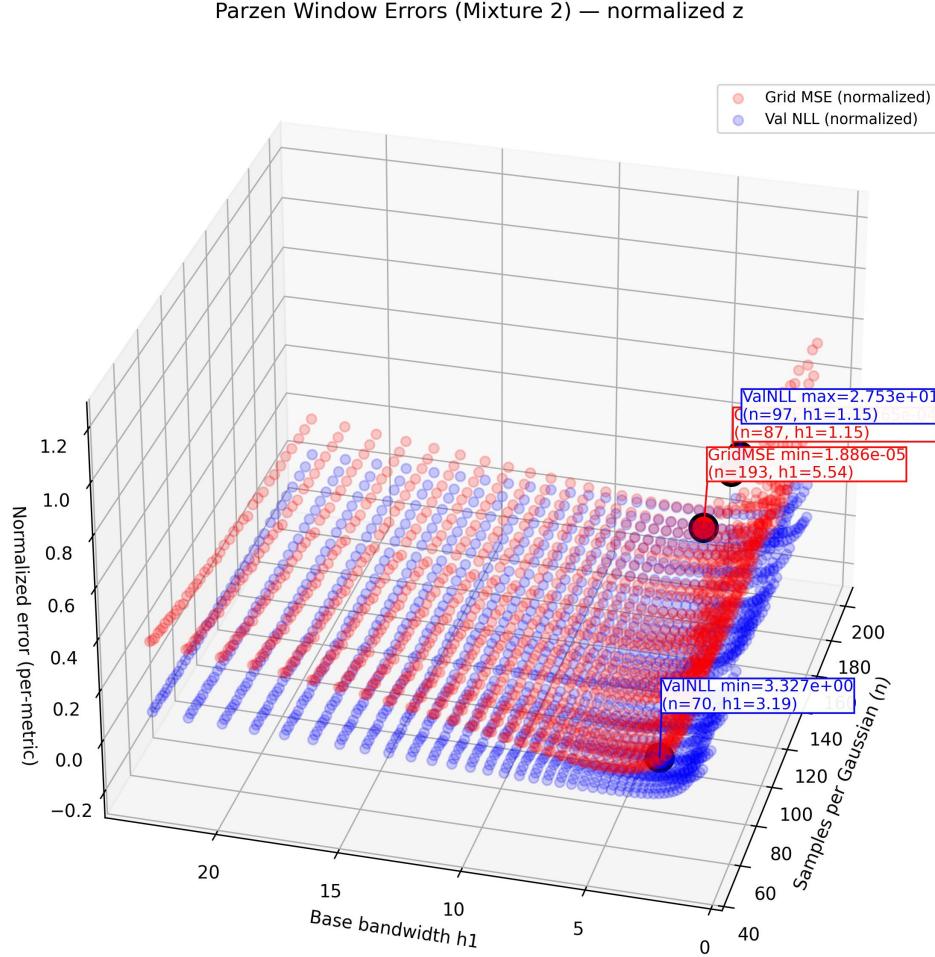


Figure 2: In red: MSE between the mixture 2 pdf and its estimate PW estimate  
 In blue: ValNLL between PW estimate and the sampled points;  
 while varying the window size  $h_1$  and the number of sampled points for each gaussian.

In this graph it's also noticeable how the MSE steepnes in the oversmoothed region is less pronounced compared to mixture 1; this is probably due to the fact that oversmoothing the probability mass causes less error when the modes are closer together; whilst, in the mixture 1, oversmoothing causes the probability mass to fall on the tails; which generates more error.

### Parzen Window Errors (Mixture 3) — normalized z

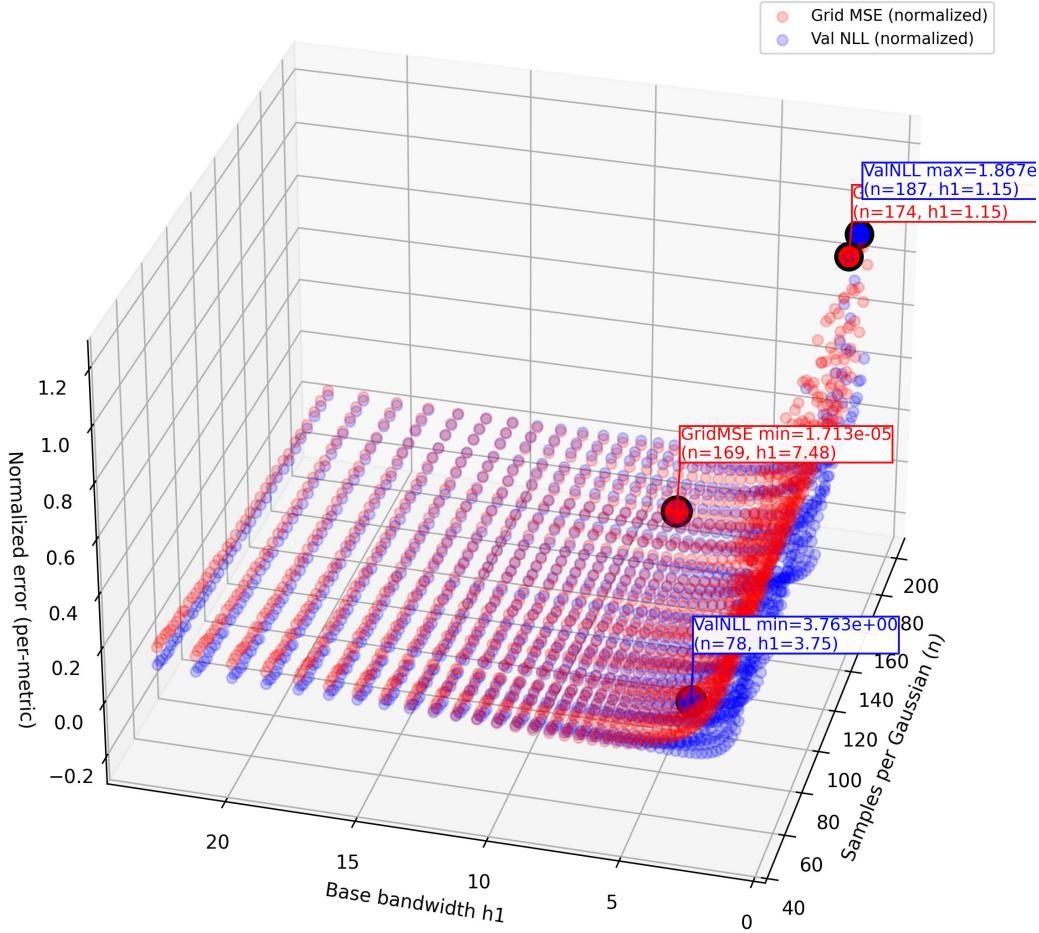


Figure 3: In red: MSE between the mixture 3 pdf and its estimate PW estimate  
 In blue: ValNLL between PW estimate and the sampled points;  
 while varying the window size  $h_1$  and the number of sampled points for each gaussian.

In this graph it's see how increasing sampled points seems to have less impact on reducing MSE compared to mixture 1 and 2. NLL does also not seem to change neither the shape in the undersmoothing region nor the optimal  $h_1$  value. whilst the real optimal  $h_1$  for the MSE seems to change with the number of gaussians in the mixture. The NLL does not, causing the selected mixtures to be more oversmoothed, the more gaussians there are in the mixture.

$$h_{1,\text{MSE}}^{\text{opt}}(\text{mix}_1, \text{mix}_2, \text{mix}_3) = (5.13, 5.54, 7, 48), \quad h_{1,\text{NLL}}^{\text{opt}}(\text{mix}_1, \text{mix}_2, \text{mix}_3) = (2.11, 3.19, 3.75).$$

### 1.1.2 Parzen Window Overlays

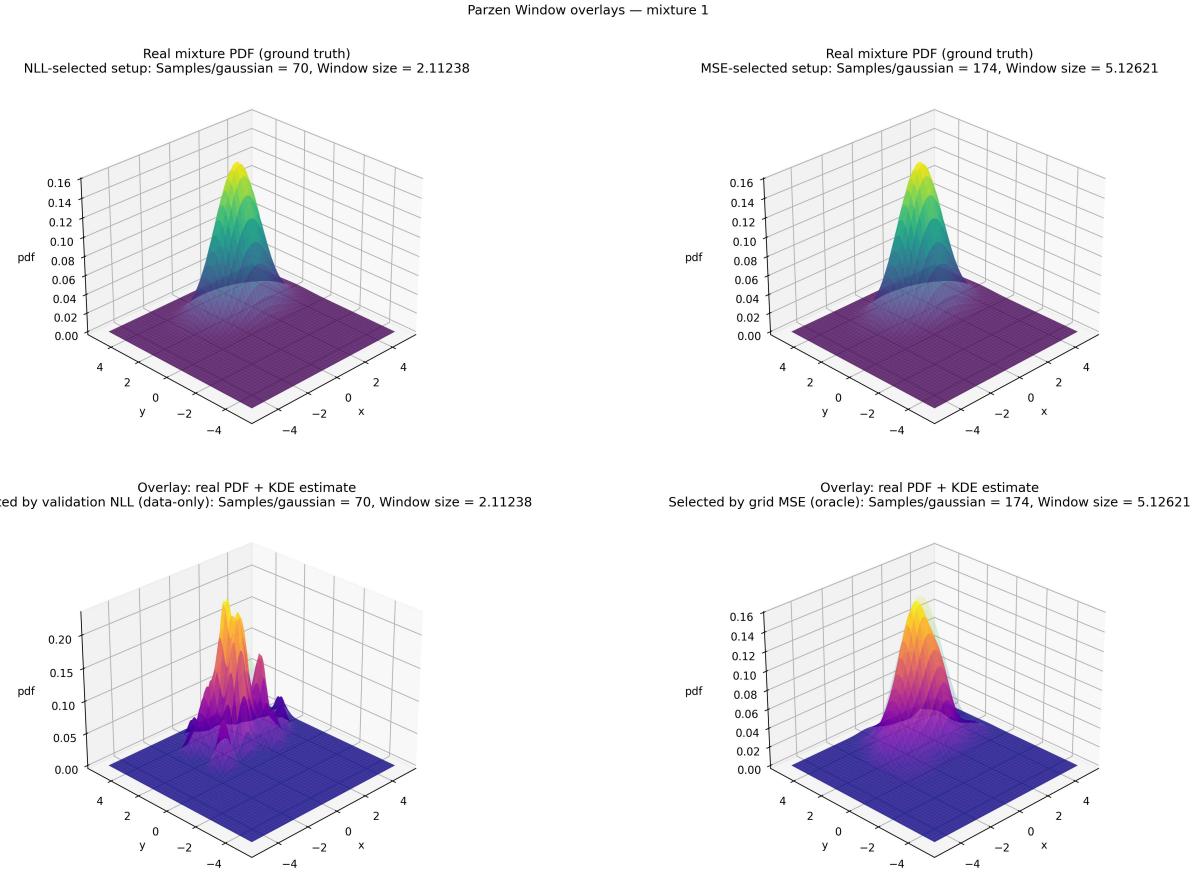


Figure 4: Figure displaying the overlay between the real pdf of mixture 1 and its PW estimate; with the optimal parameters selected by MSE and NLL respectively.

In this graph it's noticeable how increasing the sampled points per gaussian does help in reducing the MSE; but the effect on the NLL is less visible.  
Therefore the selected parameters by NLL are undersmoothed compared to the MSE ones.

Parzen Window overlays — mixture 2

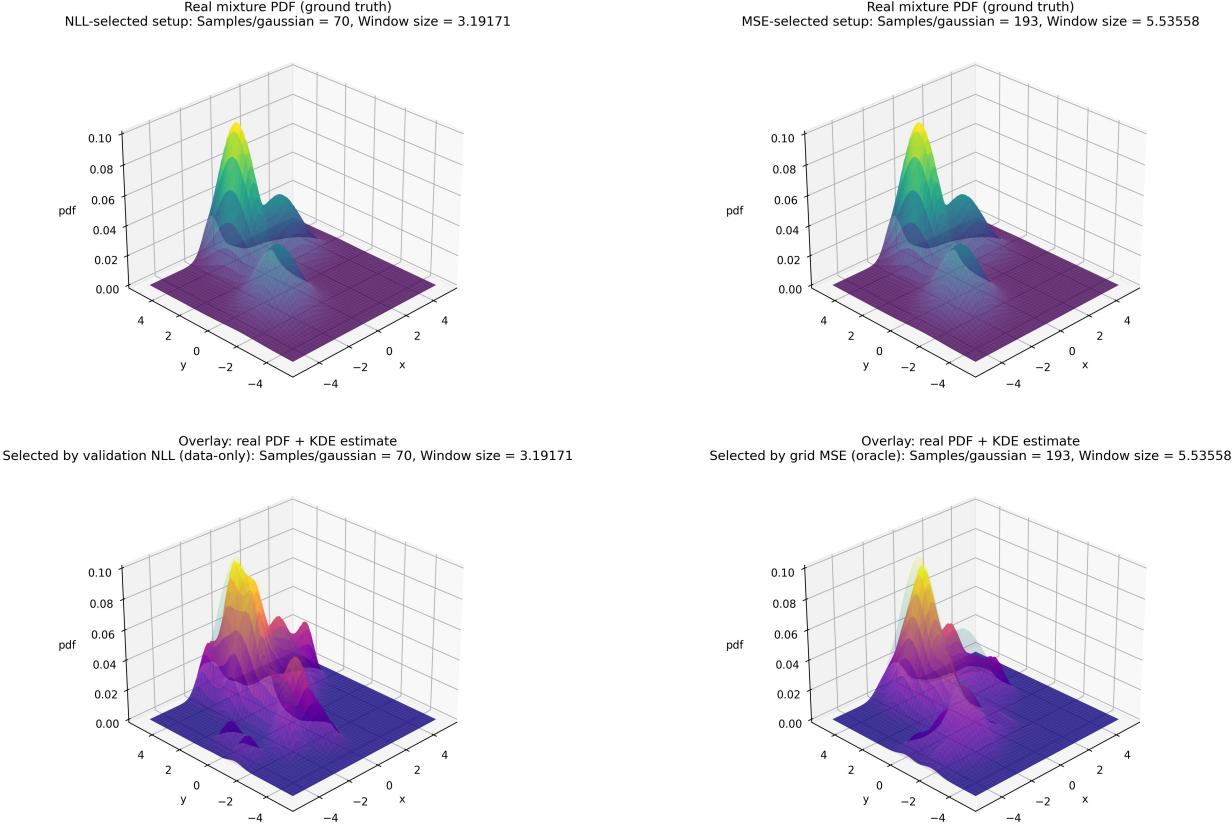


Figure 5: Figure displaying the overlay between the real pdf of mixture 2 and its PW estimate; with the optimal parameters selected by MSE and NLL respectively.

It's also noticeable how the MSE selected parameters still cannot provide high accuracy in estimating the pdf's gaussians with very different variances, this is due to the fact that  $h_1$  is constant across the whole input space. therefore high peaks result oversmoothed to accommodate for the wider modes and viceversa.

Parzen Window overlays — mixture 3

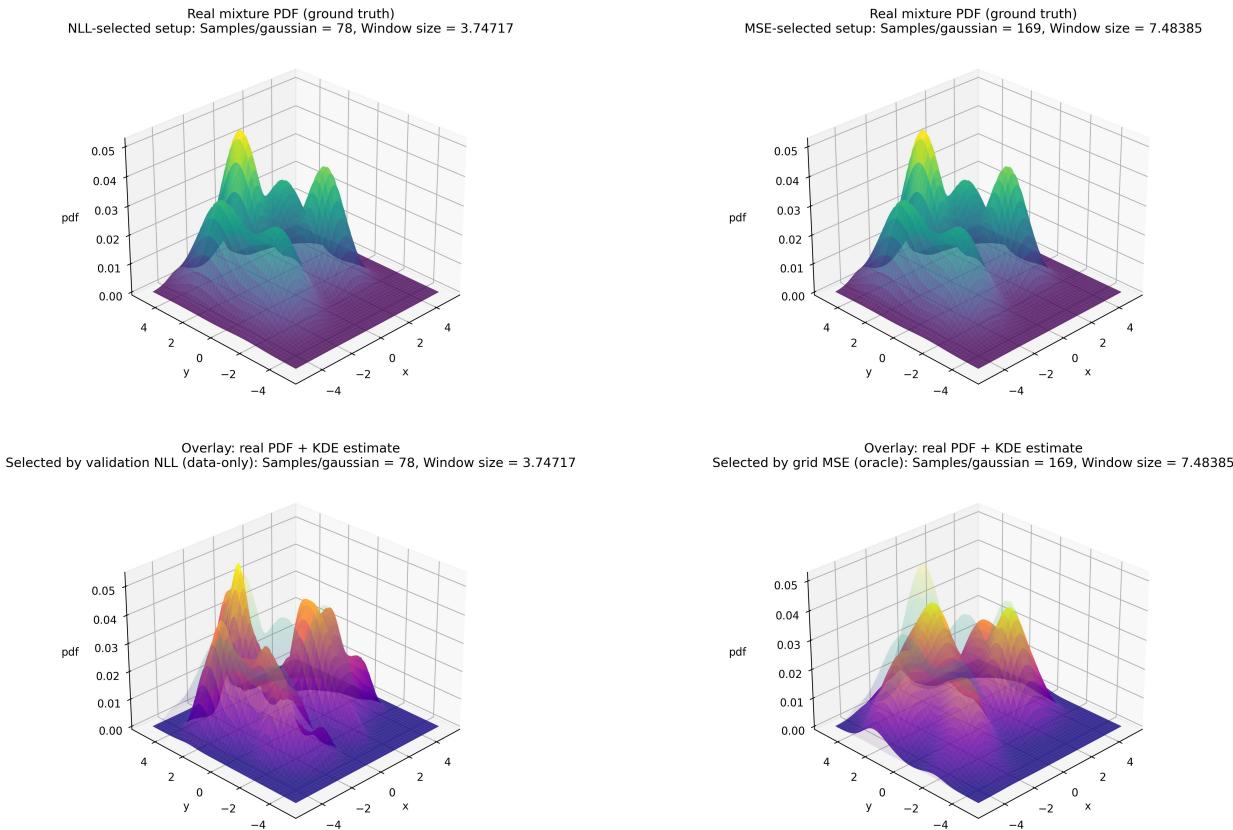


Figure 6: Figure displaying the overlay between the real pdf of mixture 2 and its PW estimate; with the optimal parameters selected by MSE and NLL respectively.

In this graph it's instead noticeable how NLL selected parameters may sometime lead the predicted pdf with undersmoothed regions; that approximates single peaks as two distinct modes.

This could even be whilst having a low enough NLL value; this is because data points drawn from those peaks will still have high likelihood even if there is a valley in between them.

This effect is less visible in MSE selected parameters since the overall shape of the pdf is more important than the likelihood of single data points.

## 1.2 Parzen Neural Network

### 1.2.1 Parzen Neural Network Errors

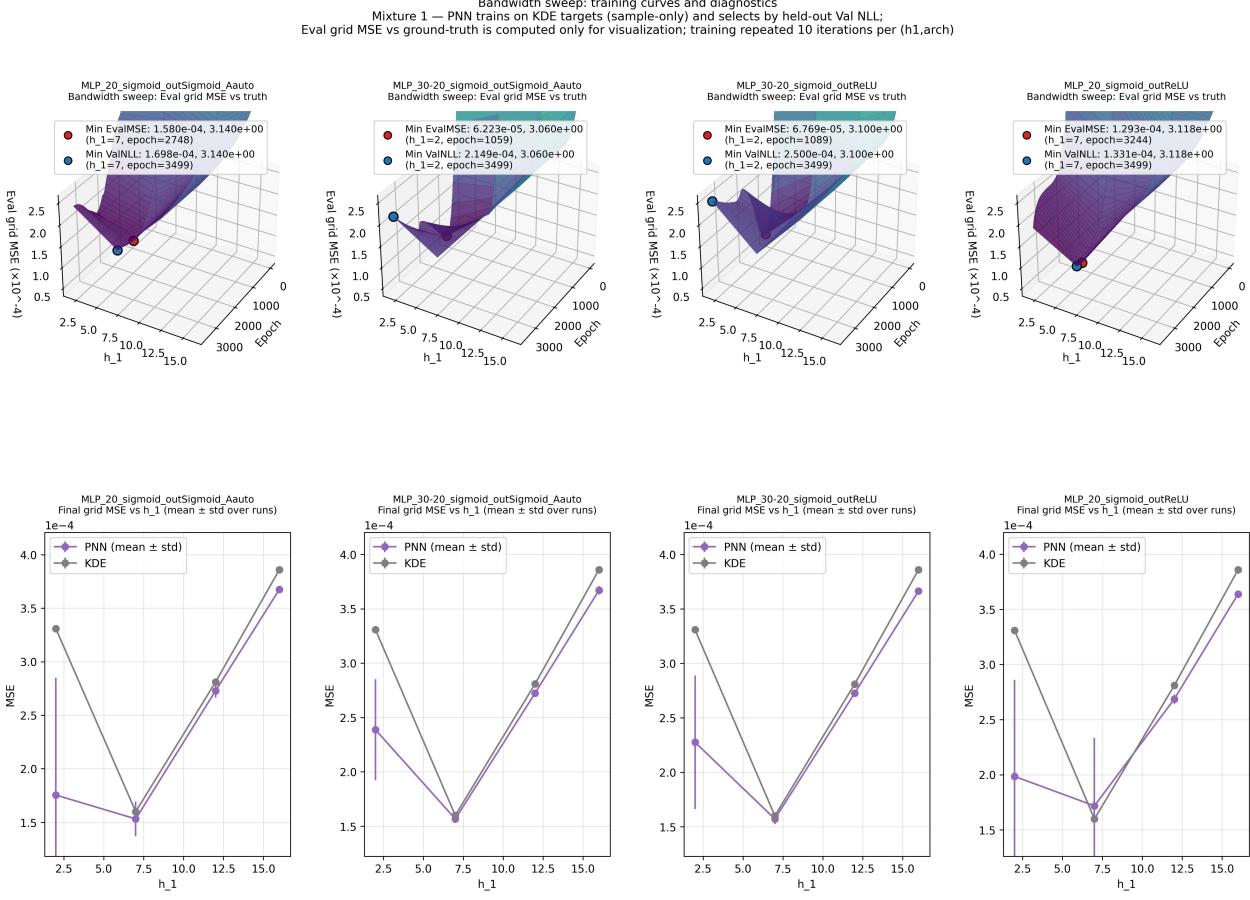


Figure 7: Top: Eval-MSE surface between mixture 1 pdf and PNN estimate over epoch  $\times h_1$   
Bottom: final grid MSE (last epoch) vs  $h_1$  shown as mean  $\pm$  std  
across 10 runs (PNN) and KDE for reference.

In this graphs it's noticeable how the PNNs tends to have a lower optimal bandwidth with respect to the PW.  
which could also be seen in the error subgraphs with MSE vs  $h_1$ ,  
where it is obvious that the PNN seems to work better with undersmoothed KDEs.

It is also marked the difference between the std of the MSE at a certain  $h_1$  with deep architectures,  
where the result between iterations, with the same epochs and the same  $h_1$  might have very different results.

Bandwidth sweep: training curves and diagnostics  
Mixture 2 — PNN trains on KDE targets (sample-only) and selects by held-out Val NLL;  
Eval grid MSE vs ground-truth is computed only for visualization; training repeated 10 iterations per (h1,arch)

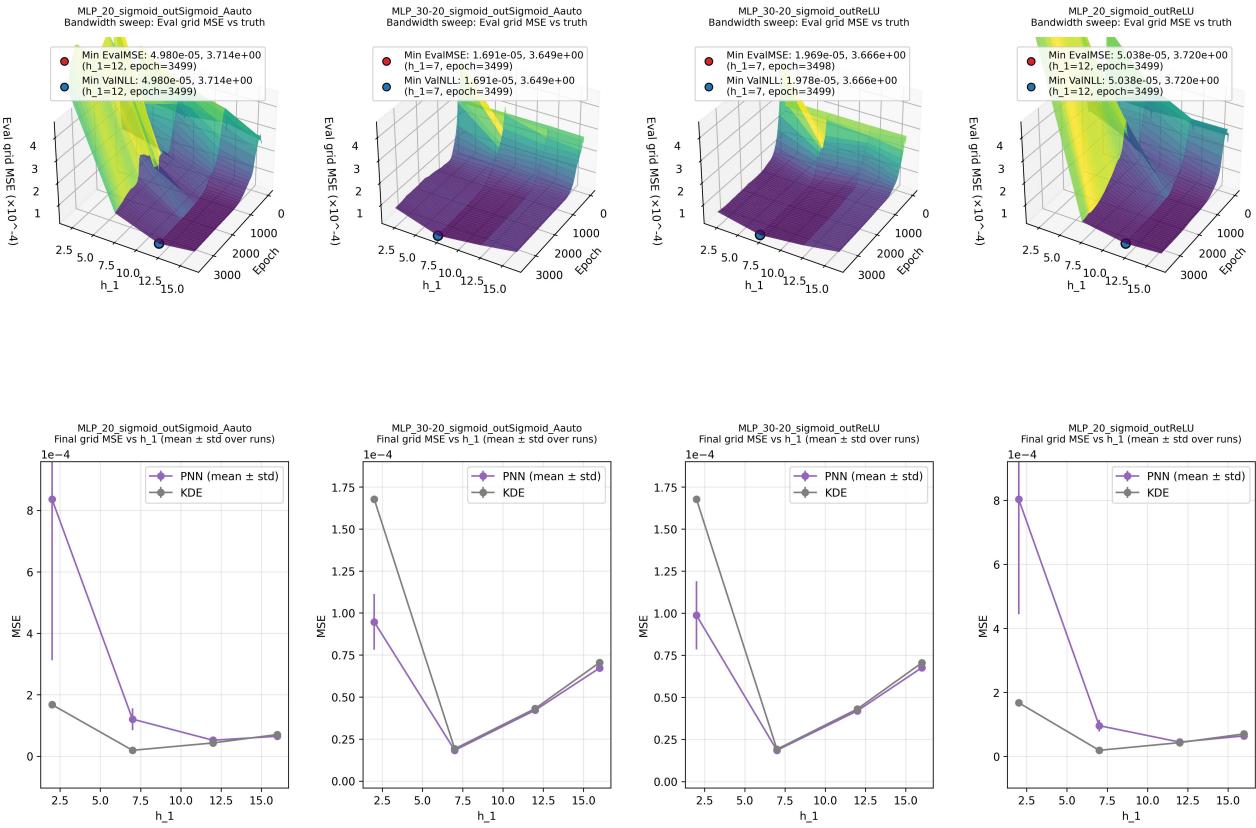


Figure 8: Top: Eval-MSE surface between mixture 2 pdf and PNN estimate over epoch  $\times h_1$   
Bottom: final grid MSE (last epoch) vs  $h_1$  shown as mean  $\pm$  std  
across 10 runs (PNN) and KDE for reference.

In this graph it's noticeable how increasing the number of gaussians in the mixture causes:

- Single hidden layer PNNs to perform worse than KDEs and double hidden layered PNNs, whilst being extremely sensitive to  $h_1$  variations, especially in the undersmoothed region.
- The MSE mesh seems to be flattened, especially in the oversmoothed region, with respect to the mixture1.
- In this graph it's marked how PNNs with deeper architectures, have a high resiliency to having undersmoothed KDEs, whilst it still causes overfitting if the stopping rule is not set.

Bandwidth sweep: training curves and diagnostics  
Mixture 3 — PNN trains on KDE targets (sample-only) and selects by held-out Val NLL;  
Eval grid MSE vs ground-truth is computed only for visualization; training repeated 10 iterations per (h1,arch)

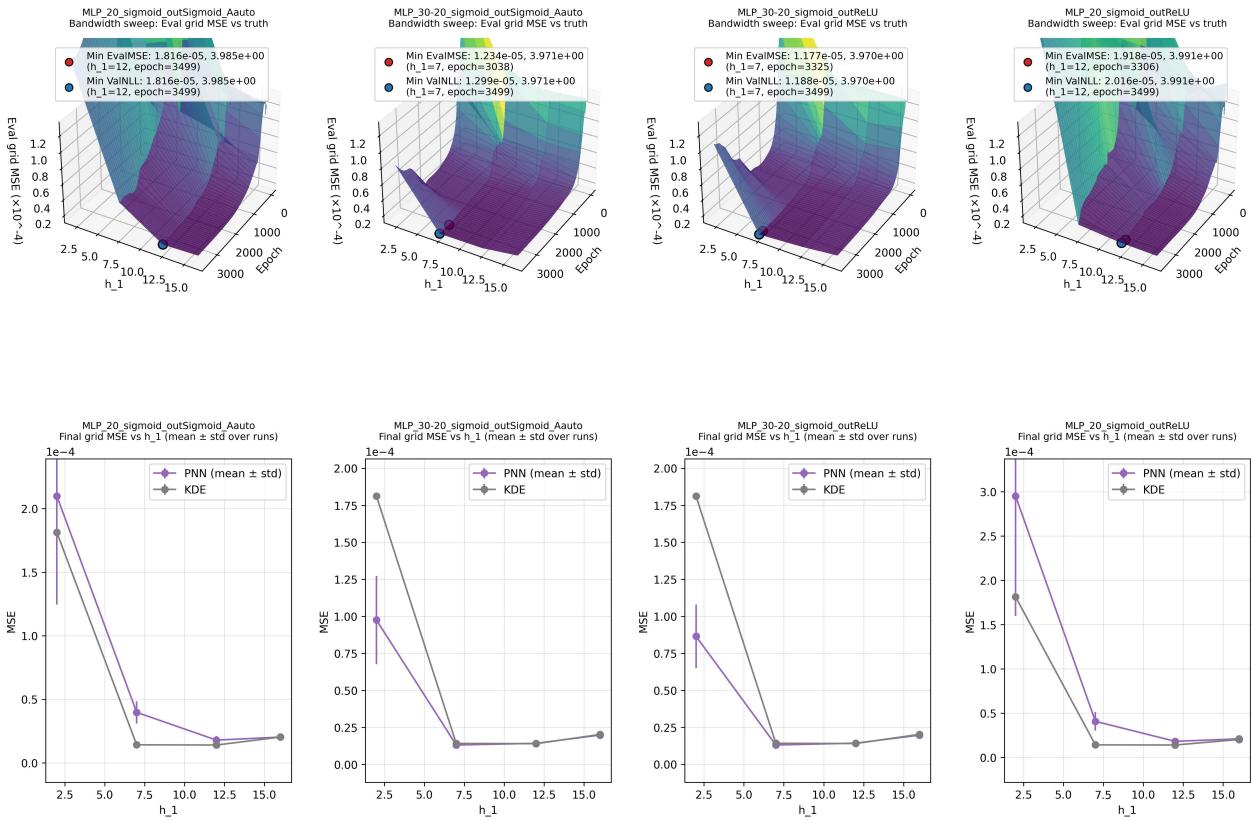


Figure 9: Top row: Eval-MSE surface between mixture 3 pdf and PNN estimate over epoch  $\times h_1$   
Bottom row: final grid MSE (last epoch) vs  $h_1$  shown as mean  $\pm$  std  
across 10 runs (PNN) and KDE for reference.

It's marked in this graphs how increasing the number of points inside the PDF (by increasing the number of gaussians inside the mixture); causes the Val-NLL (based onto held-out data points) to be a good metric to determine the best combination of architecture, stopping epoch and KDE bandwidth. This could be seen because the mixture 3 is the one where Val-NLL points have the lowest EvalMSE, across all mixtures.

### 1.2.2 Parzen Neural Network Overlays

In this subsection are shown the overlays at minimum validation NLL (ValNLL), for each mixture and each architecture.

Each column corresponds to one architecture and is displayed at its own best bandwidth  $h_1$  (the one minimizing ValNLL for that architecture on the held-out split).

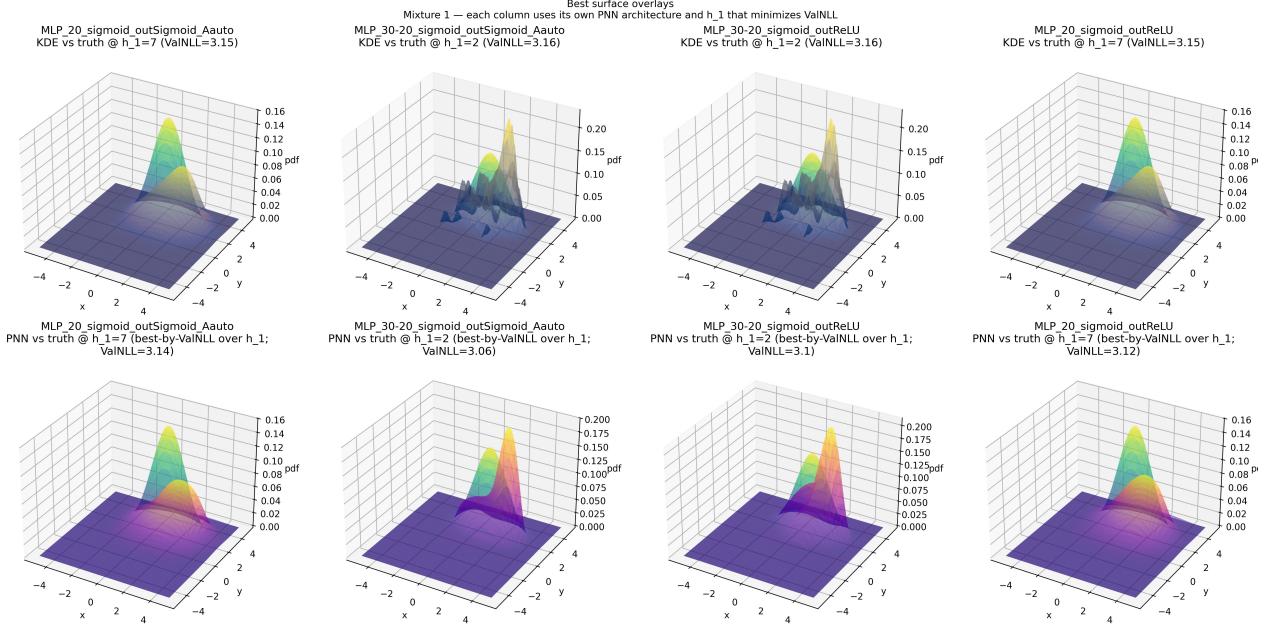


Figure 10: Best overlays for mixture 1, selected by minimizing ValNLL.

Top row: Parzen Window (KDE) estimate vs ground-truth mixture PDF at the same selected  $h_1$ .

Bottom row: Parzen Neural Network (PNN) estimate vs ground-truth mixture PDF.

In this graph it is appreciable how PNN overlays from deep architectures have an easier time at approximating high variance sections.

and how using held-out points for validation causes selection of smoothed out overlays even if target KDEs are undersmoothed.

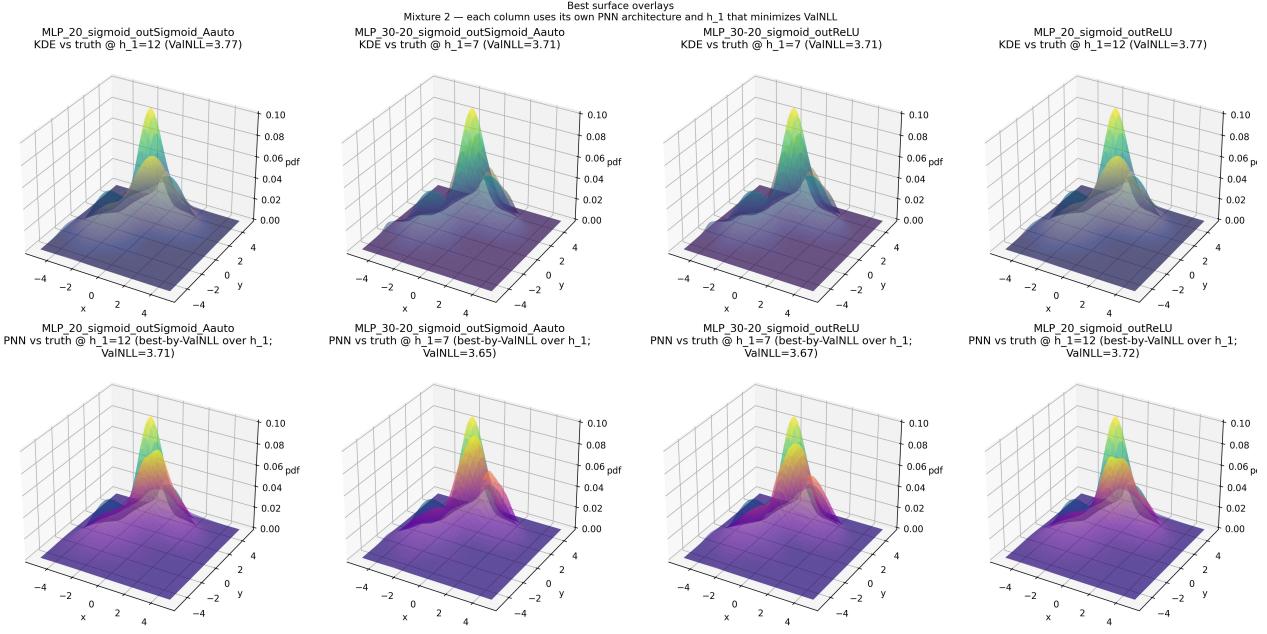


Figure 11: Best overlays for mixture 2, selected by minimizing ValNLL.

Top row: Parzen Window (KDE) estimate vs ground-truth mixture PDF at the same selected  $h_1$ .  
 Bottom row: Parzen Neural Network (PNN) estimate vs ground-truth mixture PDF.

In this graphs it's instead marked how single layered PNNs perform better with oversmoothed KDEs, even if this causes poor MSE with the ground truth.????

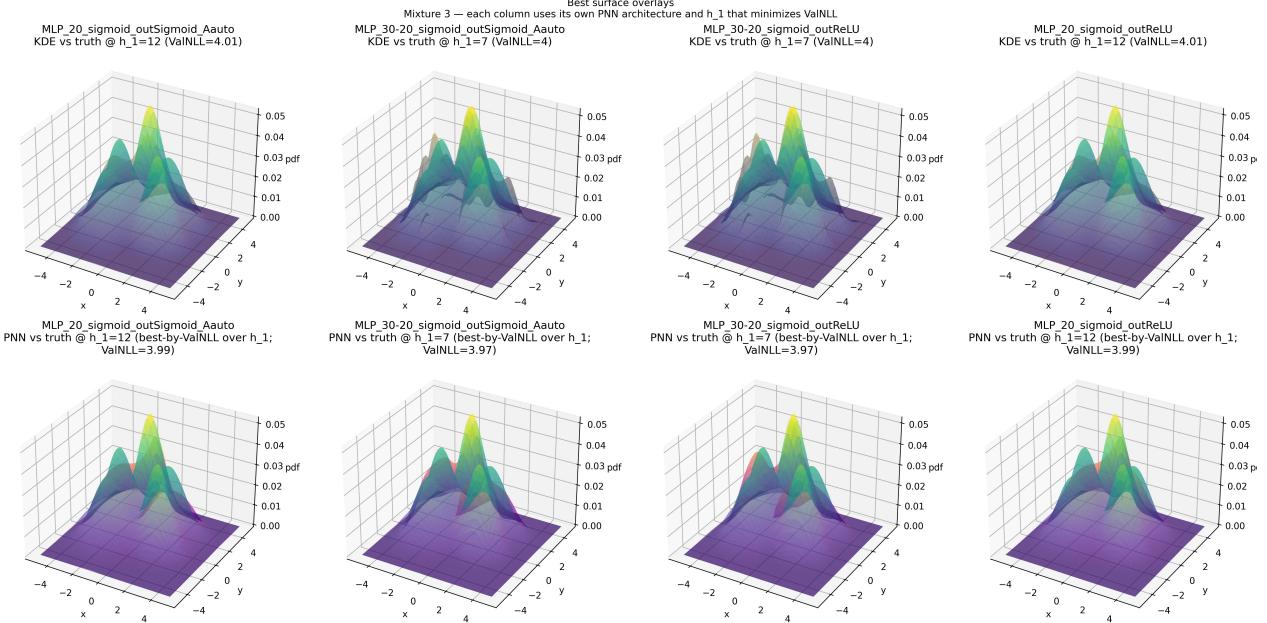


Figure 12: Best overlays for mixture 3, selected by minimizing ValNLL.

Top row: Parzen Window (KDE) estimate vs ground-truth mixture PDF at the same selected  $h_1$ .  
 Bottom row: Parzen Neural Network (PNN) estimate vs ground-truth mixture PDF.

It's marked how in this graphs the PNNs with higher layer count still have a hard time at approximating high variance peaks, close to each others in the pdf, this could be due to the fact that giving that most points are inside the  $R_n$  region inside the peak, the MLE

that estimates points between the peaks, as part of the peaks, still get overall low NLL because those points are few, compared to the ones inside the peaks