# Supplemental Material for Active Decision Boundary Annotation with Deep Generative Models

Miriam Huijser
Aiir Innovations
Amsterdam, The Netherlands
https://aiir.nl/

Jan C. van Gemert
Delft University of Technology
Delft, The Netherlands
http://jvgemert.github.io/

## 1. Projecting query point on decision boundary

Projecting query point $\mathbf{z}^*$ on the decision boundary parameterized by $\mathbf{w}$ and $b$ yields $\mathbf{z}^p$, see figure 1.
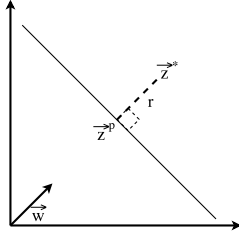


Figure 1: $\mathbf{z}^p$ is the projection of query point $\mathbf{z}^*$ on the decision boundary.

The projected point is defined as

$$\mathbf{z}^p = \mathbf{z}^* - yr\frac{\mathbf{w}}{|\mathbf{w}|}, \tag{1}$$

where $y$ is the class label of $\mathbf{z}^*$. Furthermore, $\mathbf{z}^p$ lies on the decision boundary and thus satisfies

$$\mathbf{w}^\intercal \mathbf{z}^p + b = 0. \tag{2}$$

Substituting Eq. (1) in Eq. (2) yields:

$$\mathbf{w}^\intercal \left( \mathbf{z}^* - yr\frac{\mathbf{w}}{|\mathbf{w}|} \right) + b = 0. \tag{3}$$

Now, we solve for $r$:

$$r = \frac{\mathbf{w}^\intercal \mathbf{z}^* + b}{y|\mathbf{w}|}. \tag{4}$$

Substituting $r$ from Eq. (4) in Eq. (1) gives:

$$\begin{aligned} \mathbf{z}^p &= \mathbf{z}^* - y\frac{\mathbf{w}^\intercal \mathbf{z}^* + b}{y|\mathbf{w}|} \cdot \frac{\mathbf{w}}{|\mathbf{w}|} \\ &= \mathbf{z}^* - \frac{(\mathbf{w}^\intercal \mathbf{z}^* + b)\mathbf{w}}{|\mathbf{w}||\mathbf{w}|} \\ &= \mathbf{z}^* - \frac{(\mathbf{w}^\intercal \mathbf{z}^* + b)\mathbf{w}}{\mathbf{w}^\intercal \mathbf{w}}. \end{aligned} \tag{5}$$

## 2. Results measured in Average Precision

As requested by one of our anonymous reviewer we also evaluate our results measured in Average Precision, see Tables 1, 2, 3, 4.

| Experiment 5: Full dataset evaluation | | |
|---|---|---|
| | Sample | Boundary (ours) |
| MNIST | $99.78 \pm 0.02$ | $\mathbf{99.84 \pm 0.03}$ |
| SVHN | $90.65 \pm 0.19$ | $91.22 \pm 1.95$ |
| Shoe-Bag | $98.89 \pm 0.18$ | $\mathbf{99.34 \pm 0.13}$ |

Table 4: Average Precision results for sample-based active learning and boundary active learning for all datasets averaged over 150 queries (maximum possible score is 150), averaged over all class pairs. The experiments are repeated 5 times and significant results are shown in bold. Significance is measured with a paired t-test with p $< 0.05$.

Experiment 1: Evaluating various query strategies (average precision)

| Strategy | MNIST 0 vs. 8 | | SVHN 0 vs. 8 | | Shoe-Bag | |
| | Sample | Boundary (ours) | Sample | Boundary (ours) | Sample | Boundary (ours) |
| --- | --- | --- | --- | --- | --- | --- |
| Uncertainty | $98.9 \pm 0.2$ | $\mathbf{99.4 \pm 0.2}$ | $91.1 \pm 0.8$ | $\mathbf{93.8 \pm 0.7}$ | $98.9 \pm 0.2$ | $\mathbf{99.3 \pm 0.1}$ |
| Uncertainty-dense | $94.5 \pm 11.0$ | $\mathbf{96.5 \pm 10.7}$ | $77.1 \pm 5.7$ | $\mathbf{89.4 \pm 2.1}$ | $84.6 \pm 4.0$ | $\mathbf{96.9 \pm 1.1}$ |
| 5 Cluster centroid | $98.6 \pm 0.1$ | $\mathbf{99.6 \pm 0.03}$ | $75.5 \pm 5.4$ | $\mathbf{83.1 \pm 1.2}$ | $95.1 \pm 0.8$ | $\mathbf{99.3 \pm 0.1}$ |
| Random | $98.4 \pm 0.4$ | $\mathbf{99.2 \pm 0.3}$ | $89.3 \pm 1.4$ | $\mathbf{93.7 \pm 0.8}$ | $98.1 \pm 0.4$ | $\mathbf{99.3 \pm 0.1}$ |

Table 1: Average precision averaged over 150 iterations.

Experiment 3: Evaluating annotation noise

| Sampling noise (# images) | MNIST 0 vs. 8 | | SVHN 0 vs. 8 | | Shoe-Bag | |
| | Sample | Boundary (ours) | Sample | Boundary (ours) | Sample | Boundary (ours) |
| --- | --- | --- | --- | --- | --- | --- |
| 0 | $99.0 \pm 0.2$ | $\mathbf{99.5 \pm 0.1}$ | $90.9 \pm 1.4$ | $\mathbf{93.7 \pm 0.8}$ | $98.9 \pm 0.2$ | $\mathbf{99.3 \pm 0.2}$ |
| 1 | $99.0 \pm 0.2$ | $\mathbf{99.5 \pm 0.1}$ | $90.9 \pm 1.4$ | $\mathbf{93.4 \pm 0.7}$ | $98.9 \pm 0.2$ | $\mathbf{99.3 \pm 0.2}$ |
| 2 | $99.0 \pm 0.2$ | $\mathbf{99.4 \pm 0.2}$ | $90.9 \pm 1.4$ | $\mathbf{92.3 \pm 1.6}$ | $98.9 \pm 0.2$ | $99.1 \pm 0.6$ |
| 3 | $99.0 \pm 0.2$ | $\mathbf{99.3 \pm 0.1}$ | $90.9 \pm 1.4$ | $\mathbf{92.5 \pm 0.8}$ | $98.9 \pm 0.2$ | $99.1 \pm 0.6$ |
| 4 | $99.0 \pm 0.2$ | $\mathbf{99.1 \pm 0.1}$ | $90.9 \pm 1.4$ | $91.3 \pm 0.9$ | $98.9 \pm 0.2$ | $\mathbf{99.1 \pm 0.3}$ |
| 5 | $99.0 \pm 0.2$ | $99.0 \pm 0.1$ | $90.9 \pm 1.4$ | $86.7 \pm 9.7$ | $98.9 \pm 0.2$ | $98.8 \pm 0.4$ |

Table 2: Average Precision results for noisy boundary active learning with uncertainty sampling for MNIST (classifying 0 and 8), SVHN (classifying 0 and 8) and Handbags vs. Shoes averaged over 150 queries (maximum possible score is 150). Each experiment is repeated 15 times. For each row, the significantly best result is shown in bold, where significance is measured with a paired t-test with $p < 0.05$. Noise has been added to the boundary annotation points; not to the image labels.

Experiment 4: Evaluating a human oracle

| Annotation | MNIST 0 vs. 8 | | SVHN 0 vs. 8 | | Shoe-Bag | |
| | Sample | Boundary (ours) | Sample | Boundary (ours) | Sample | Boundary (ours) |
| --- | --- | --- | --- | --- | --- | --- |
| Human oracle | - | - | $63.0 \pm 7.2$ | $64.6 \pm 7.6$ | - | - |
| SVM oracle | $94.7 \pm 2.5$ | $\mathbf{96.1 \pm 2.7}$ | $74.0 \pm 5.8$ | $74.9 \pm 5.4$ | $94.9 \pm 1.5$ | $\mathbf{95.8 \pm 1.6}$ |

Table 3: Average Precision results for a human and a SVM oracle for sample-based active learning and our boundary active learning for MNIST (classifying 0 and 8), SVHN (classifying 0 and 8) and Shoe-Bag averaged over 10 queries (maximum possible score is 10). The experiments are repeated 15 times and significant results per row are shown in bold for $p < 0.05$.