

Bergische Universität Wuppertal

Fachbereich Mathematik und Naturwissenschaften

Institute of Mathematical Modelling, Analysis and Computational
Mathematics (IMACM)

Preprint BUW-IMACM 18/18

Matthias Rottmann, Pascal Colling, Thomas Paul Hack, Fabian
Hüger, Peter Schlicht and Hanno Gottschalk

**Prediction Error Meta Classification in Semantic
Segmentation: Detection via Aggregated
Dispersion Measures of Softmax Probabilities**

December 21, 2018

<http://www.math.uni-wuppertal.de>

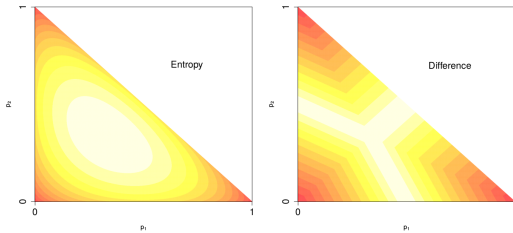


Figure 1: Visualization of entropy (cf. (2), left panel) and difference (cf. (3), right hand panel) for three variables p_1 , p_2 and p_3 treating p_3 implicitly via $p_1 + p_2 + p_3 = 1$.

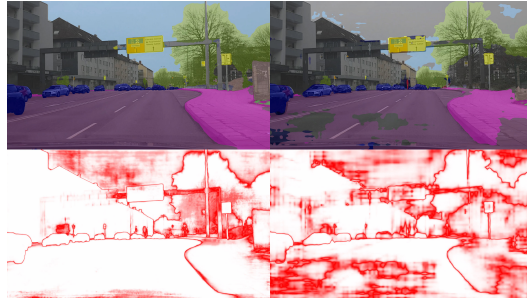


Figure 2: Example of segmentation (top line) and heat map D_z (bottom line) for Xception65 (left column) and MobilenetV2 (right hand column). Original image is not part of the Cityscapes dataset.

the context of classical machine learning for learning the weights for each member of a committee of classifiers [14]. In terms of deep learning we use it as a shorthand to distinguish between a network’s own classification and the classification whether a prediction is “true” or “false”. In contrast to the work cited above, we aim at judging the statistical reliability of each segment inferred by the neural network. To this end, dispersion measures, like entropy, are applied to the softmax probabilities (the network’s output) on pixel level yielding dispersion heat maps. We aggregate these heat maps over predicted segments alongside with other quantities derived from the network’s prediction like the segment’s size and predicted class. From this, we construct per-segment metrics. A commonly used performance measure for the quality of a segmentation is the intersection over union (IoU a.k.a. Jaccard index [10]) of prediction and ground truth. We use the constructed metrics as inputs to logistic regression models for meta classifying, whether an inferred segment’s IoU vanishes or not, i.e., predicting $IoU = 0$ or $IoU > 0$. Also, we use linear regression models for predicting a segment’s IoU directly, thus obtaining statements about the reliability of the network’s prediction.

In our tests, we employ two different publicly available DeepLabv3+ networks [2] that were trained on the Cityscapes dataset [4]. We perform all tests on the Cityscapes validation set and demonstrate that our segment-wise metrics are well correlated with the IoU ; thus they are suitable for detecting false positives on segment level. For logistic regression fits we obtain values of up to 87.71% for the area under curve corresponding to the receiver operator characteristic curve (AUROC, see [6]). Predicting the segment-wise IoU via linear regression we obtain prediction standard deviations of down to 0.130 and R^2 values of up to 81.48%.

2 Pixel-wise dispersion metrics and aggregation over segments

A segmentation network with a softmax output layer can be seen as a statistical model that provides for each pixel z of the image a probability distribution $f_z(y|x, w)$ on the q class labels $y \in \mathcal{C} = \{y_1, \dots, y_q\}$, given the weights w and the data x . The predicted class in z is then given by

$$\hat{y}_z(x, w) = \arg \max_{y \in \mathcal{C}} f_z(y|x, w). \quad (1)$$

Dispersion or concentration measures quantify the degree of randomness in $f_z(y|x, w)$. Here, we consider two of those measures: *entropy* E_z (also known as *Shannon information* [17]) and *difference in probability* D_z , i.e. the difference between the two largest softmax values:

$$E_z(x, w) = -\frac{1}{\log(q)} \sum_{y \in \mathcal{C}} f_z(y|x, w) \log f_z(y|x, w), \quad (2)$$

$$D_z(x, w) = 1 - f_z(\hat{y}_z(x, w)|x, w) + \max_{y \in \mathcal{C} \setminus \{\hat{y}_z(x, w)\}} f_z(y|x, w). \quad (3)$$

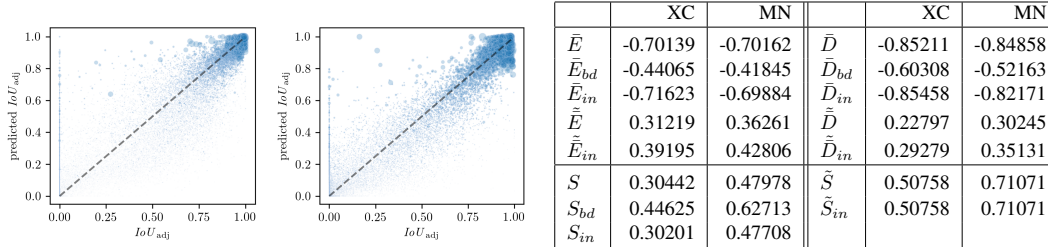


Figure 3: IoU_{adj} vs. IoU_{adj} for all connected components predicted by Xception65 (left) and MobilenetV2 (right). Dot sizes are proportional to S .

Table 1: Correlation coefficients ρ with respect to IoU_{adj} . Results are computed on the Cityscapes validation set, XC: DeepLabv3+Xception65 and MN: DeepLabv3+MobilenetV2.

For better comparison, both quantities have been written as dispersion measures and been normalized to the interval $[0, 1]$: One has $E_z = D_z = 1$ for the equiprobability distribution $f_z(y|x, w) = \frac{1}{q}$, $y \in \mathcal{C}$, and $E_z = D_z = 0$ on the deterministic probability distribution ($f_z(y|x, w) = 1$ for one class and 0 otherwise). For the discussion of further dispersion measures, cf. [5]. Figure 1 displays these quantities for three class probability distributions. The most direct method of uncertainty quantification on an image is the heat mapping of a dispersion measure as in fig. 2.

For a given image x we denote by $\hat{\mathcal{K}}_x$ the set of connected components (segments) in the predicted segmentation $\hat{\mathcal{S}}_x = \{\hat{y}_z(x, w) | z \in x\}$ (omitting the dependence on the weights w). Analogously we denote by \mathcal{K}_x the set of connected components in the ground truth \mathcal{S}_x . For each $k \in \hat{\mathcal{K}}_x$, we define the following quantities:

- the interior $k_{in} \subset k$ where a pixel z is an element of k_{in} if all eight neighbouring pixels are an element of k
- the boundary $k_{bd} = k \setminus k_{in}$
- the intersection over union IoU : let $\mathcal{K}_x|_k$ be the set of all $k' \in \mathcal{K}_x$ that have non-trivial intersection with k and whose class label equals the predicted class for k , then

$$IoU(k) = \frac{|k \cap K'|}{|k \cup K'|}, \quad K' = \bigcup_{k' \in \mathcal{K}_x|_k} k'$$

- adjusted IoU_{adj} : let $Q = \{q \in \hat{\mathcal{K}}_x : q \cap K' \neq \emptyset\}$, for reasons explained in the appendix we use in our tests

$$IoU_{adj}(k) = \frac{|k \cap K'|}{|k \cup (K' \setminus Q)|}$$

- the pixel sizes $S = |k|$, $S_{in} = |k_{in}|$, $S_{bd} = |k_{bd}|$
- the mean entropies \bar{E} , \bar{E}_{in} , \bar{E}_{bd} defined as

$$\bar{E}_{\#}(k) = \frac{1}{S_{\#}} \sum_{z \in k_{\#}} E_z(x), \quad \# \in \{-, in, bd\}$$

- the mean distances \bar{D} , \bar{D}_{in} , \bar{D}_{bd} defined in analogy to the mean entropies
- the relative sizes $\tilde{S} = S/S_{bd}$, $\tilde{S}_{in} = S_{in}/S_{bd}$
- the relative mean entropies $\tilde{E} = \bar{E}\tilde{S}$, $\tilde{E}_{in} = \bar{E}_{in}\tilde{S}_{in}$, and relative mean distances $\tilde{D} = \bar{D}\tilde{S}$, $\tilde{D}_{in} = \bar{D}_{in}\tilde{S}_{in}$

Typically, E_z and D_z are large for $z \in k_{bd}$. This motivates the separate treatment of interior and boundary measures. With the exception of IoU and IoU_{adj} , all scalar quantities defined above can be computed without the knowledge of the ground truth. Our aim is to analyze to which extent they are able to predict IoU_{adj} .

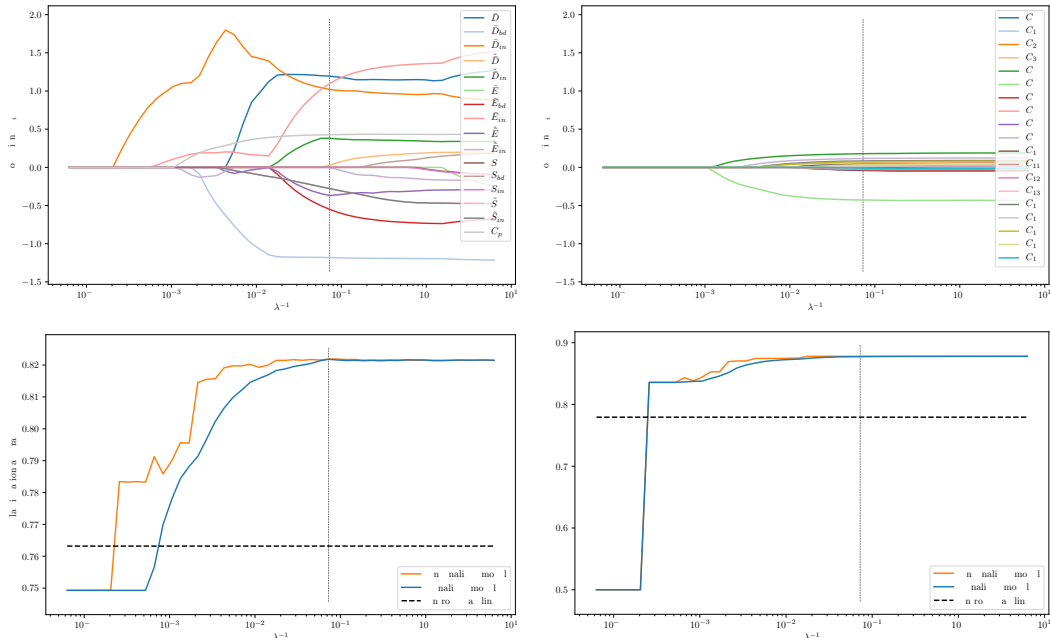


Figure 4: Results for the meta classification task $IoU_{adj} = 0, > 0$ for predictions obtained from the Xception65 net. (Top left): the weights coefficients for the 15 metrics computed with LASSO fits as function of λ^{-1} , C_p denotes the maximum of the absolute values of all weight coefficients for predicted classes. (Top right): like top left but showing coefficients for the 18 predicted classes. (Bottom left): meta classification rates for $IoU_{adj} = 0, > 0$. The blue line are the LASSO fits for different λ values, the orange line shows the performance of regular logistic regression fits ($\lambda = 0$) where the input metrics are only those that have non-zero coefficients in the LASSO fit for the current λ . (Bottom right) same as bottom left, but for AUROC. The vertical dashed lines indicate the λ value for which we obtained the best validation accuracy.

3 Numerical Experiments: Street Scenes

We investigate the properties of the metrics defined in the previous section for the example of a semantic segmentation of street scenes. To this end, we consider the DeepLabv3+ network [2] for which we use a reference implementation in Tensorflow [1] as well as weights pretrained on the Cityscapes dataset [4] and available on GitHub. The DeepLabv3+ implementation and weights are available for two network backbones: Xception65, which is a modified version of Xception [3] and is a powerful structure intended for server-side deployment, and MobilenetV2 [16], a fast structure designed for mobile devices. Each of these implementations have parameters tuning the segmentation accuracy. We choose the following best (for Xception65) and worst (for MobilenetV2) parameters in order to perform our analysis on two very distinct networks. Note, that the parameter set for the Xception65 setting also includes the evaluation of the input on multiple scales (averaging the results) which increases the accuracy and also leverages classification uncertainty. We refer to [2] for a detailed explanation of the chosen parameters.

- DeepLabv3+Xception65: output stride 8, decoder output stride 4, evaluation on input scales 0.75, 1.00, 1.25 – $mIoU = 80.42\%$ on the Cityscapes validation set
- DeepLabv3+MobilenetV2: output stride 16, evaluation on input scale 1.00 – $mIoU = 70.71\%$ on the Cityscapes validation set

Example segmentations and heat maps of the two networks are displayed in fig. 2. For both networks, we consider the output probabilities and predictions on the Cityscapes validation set, which consists of 500 street scene images at a resolution of 2048×1024 . We compute the 15 constructed metrics as well as IoU_{adj} for each segment in the segmentations of the images. In order to investigate the predictive power of the metrics, we first compute the Pearson correlation $\rho \in [-1, 1]$ between each feature and IoU_{adj} . We report the results of this analysis in table 1 and provide scatter plots of all

features relative to IoU_{adj} in [fig. 6](#). Note, that in all computations, we only consider connected components with non-empty interior.

For both networks IoU_{adj} shows strong correlation with the mean distances \bar{D} and \bar{D}_{in} as well as with the mean entropies \bar{E} and \bar{E}_{in} . On the other hand, the relative counterparts are less correlated with IoU_{adj} . The relative segment size \tilde{S} for the DeepLabv3+MobilenetV2 network shows a clear correlation whereas this is not the case for the more powerful DeepLabv3+Xception65 network.

In order to find more indicative measures, we now investigate the predictive power of the metrics when they are combined. For the Xception65 net, we obtain 45194 segments with non-empty interior of which 11331 have $IoU_{\text{adj}} = 0$. For the weaker MobilenetV2 this ratio is 42261/17671. We would first like to detect segments with $IoU_{\text{adj}} = 0$, i.e., learn the meta classification task of identifying false positive segments based on our 15 metrics and the segment-wise averaged probability distribution vectors. We term these (standardized) inputs x_k for a segment k . Further, let $y_k = \text{ceil}(IoU_{\text{adj}}) = \{0 \text{ if } IoU_{\text{adj}} = 0, 1 \text{ if } IoU_{\text{adj}} > 0\}$. The least absolute shrinkage and selection operator (LASSO, [\[18\]](#)) is a popular tool for investigating the predictive power of different combinations of input variables. We compute a series of LASSO fits, i.e., ℓ_1 -penalized logistic regression fits

$$\min_w \left[\sum_i -y_i \log(\tau(w^T x_i)) - (1 - y_i)(1 - \log(\tau(w^T x_i))) + \lambda \|w\|_1 \right], \quad (4)$$

for different regularization parameters λ and standardized inputs (zero mean and unit standard deviation). Here, $\tau(\cdot)$ is the logistic function. Results for the Xception65 net are shown in [fig. 4](#).

The top left and top right panels show, in which order the weight coefficients w for each metric/predicted class become active. At the same time the bottom left and bottom right panels show, which weight coefficient causes which amount of increase in predictive performance in terms of meta-classification rate and AUROC, respectively. The AUROC is obtained by varying the decision threshold of the logistic regression output for deciding whether $IoU = 0$ or $IoU > 0$.

The first non-zero coefficient activates the \bar{D}_{in} metric, which elevates the predictive power above our reference benchmark of choice, the mean entropy per component \bar{E} , which we term *entropy baseline*. Another significant gain is achieved when \bar{D}_{bd} and the predicted classes come into play. Noteworthy we obtain a meta-classification validation accuracy of up to 81.91% ($\pm 0.13\%$) and an AUROC of up to 87.71% ($\pm 0.15\%$) for Xception65. And also for the weaker MobilenetV2 we obtain 78.93% ($\pm 0.17\%$) classification accuracy and 86.77% ($\pm 0.17\%$) AUROC. We randomly choose 10 50/50 training/validation data splits and average the results, the numbers in brackets denote standard deviations of the averages.

Additionally, the bottom line of [fig. 4](#) shows that there is almost no performance loss when only incorporating some of the metrics proposed by the LASSO trajectory. For both networks the classification accuracy corresponds to a logistic regression trained with unbalanced meta-classes $IoU_{\text{adj}} = 0$ and $IoU_{\text{adj}} > 0$, i.e., we did not adjust the class weights. On average (over the 10 training/validation splits) 6851 components with vanishing IoU_{adj} are detected for Xception65 while 4480 remain undetected, for MobilenetV2 this ratio is 14976/2695. These ratios can be adjusted by varying the probability thresholds for deciding between $IoU_{\text{adj}} = 0$ and $IoU_{\text{adj}} > 0$. For this reason we state results in terms of AUROC which is independent of this threshold.

Ultimately, we want to predict IoU_{adj} values for all connected components and thus model an uncertainty measure. We now resign from regularization and use a linear regression model to predict the IoU_{adj} . [Figure 3](#) depicts the quality of a single linear regression fit for each of the two segmentation networks. For MobilenetV2 we obtain an R^2 value of 81.48% ($\pm 0.23\%$) and for Xception65 74.93% ($\pm 0.22\%$). [Figure 5](#) illustrates the constructed uncertainty measure with two showcases. Averaged results over 10 runs including standard deviations σ and previous meta classification result are summarized in [table 2](#). In all cases, the presented approach clearly outperforms the entropy baseline. The linear regression models do not overfit the data and note-worthily we obtain prediction standard deviations of down to 0.130 and almost no standard deviation for the averages. The classification accuracy and AUROC results are slightly biased towards the validation results as they correspond to the particular λ value that maximizes the validation accuracy. An additional discussion on the difference (also in performance between) IoU_{adj} and IoU can be found in the appendix.

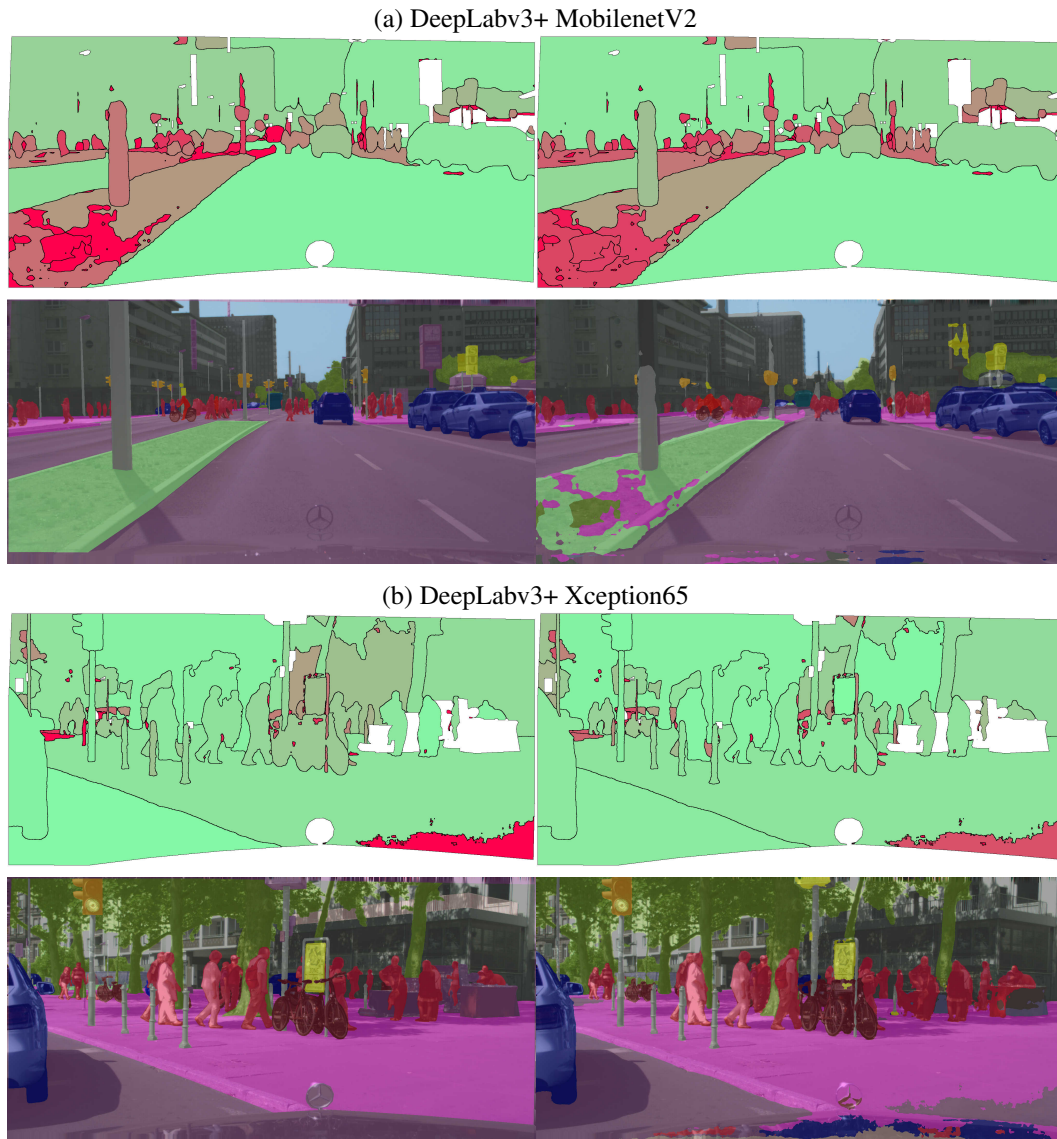


Figure 5: Prediction of the IoU_{adj} with linear regression. Each of the two sub-figures (a) and (b) consists of ground truth (bottom left), predicted segments (bottom right), true IoU_{adj} for the predicted segments (top left) and predicted IoU_{adj} for the predicted segments (top right). In the top row, green color corresponds to high IoU_{adj} values and red color to low ones, for the white regions there is no ground truth available. These regions are excluded from the statistical evaluation.

- [10] P. Jaccard. “The distribution of the flora in the alpine zone”. *New phytologist* 11.2 (Feb. 1912), pp. 37–50.
- [11] M. Kampffmeyer, A.-B. Salberg, and R. Jenssen. “Semantic segmentation of small objects and modeling of uncertainty in urban remote sensing images using deep convolutional neural networks”. *2016 ieee conference on computer vision and pattern recognition workshops (cvprw)* (2016), pp. 680–688.
- [12] A. Kendall, V. Badrinarayanan, and R. Cipolla. “Bayesian segnet: model uncertainty in deep convolutional encoder-decoder architectures for scene understanding”. *Corr abs/1511.02680* (2015).
- [13] S. Liang, Y. Li, and R. Srikant. “Principled detection of out-of-distribution examples in neural networks”. *Corr abs/1706.02690* (2017).
- [14] W. Lin and A. Hauptmann. “Meta-classification: combining multimodal classifiers”. *Mining multimedia and complex data. pakdd 2002. lecture notes in computer science 2797* (2003).

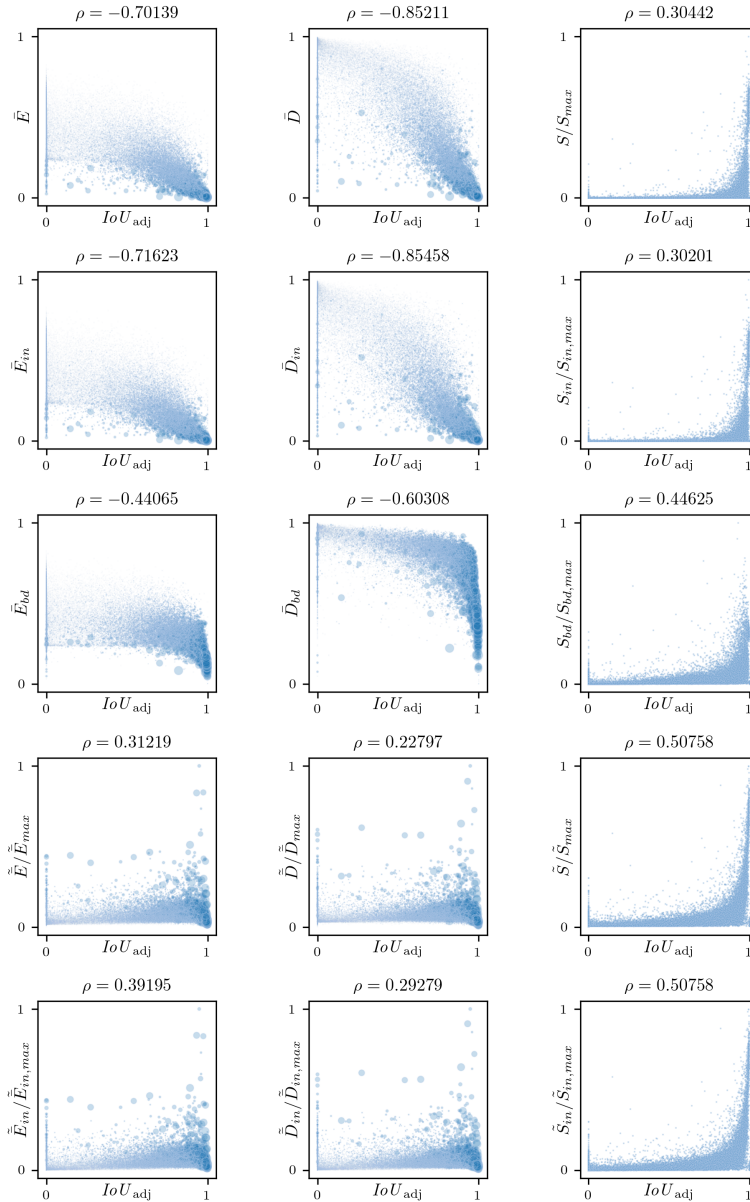


Figure 6: Correlations between IoU_{adj} and rescaled features for the DeepLabv3+Xception65 network. Dot sizes in the first two columns are proportional to S .

- [15] P. Oberdiek, M. Rottmann, and H. Gottschalk. “Classification uncertainty of deep neural networks based on gradient information”. *Artificial Neural networks and Pattern Recognition (ANNPR)*. 2018.
- [16] M. Sandler et al. “Inverted residuals and linear bottlenecks: mobile networks for classification, detection and segmentation”. *Corr abs/1801.04381* (2018).
- [17] C. E. Shannon. “A mathematical theory of communication”. In: *The Bell System Technical Journal*. Vol. 27. 1948, pp. 379–423, 623–656.
- [18] R. Tibshirani. “Regression shrinkage and selection via the lasso”. *Journal of the royal statistical society: series b* 58 (1996), pp. 267–288.
- [19] K. Wickstrøm, M. Kampffmeyer, and R. Jenssen. “Uncertainty and interpretability in convolutional neural networks for semantic segmentation of colorectal polyps”. *Corr abs/1807.10584* (2018).