# Supplementary Material

## A  Overview of Popular Regularisation Methods

First we give a brief overview of the regularisation methods tested. The main goal here is to compare a wide variety of different schemes to see if they have any common behaviour.

### A.1  Weight decay

One of the most used methods for regularisation is placing a $\ell_2$-penalty on the parameters of the network. Colloquially this is known as *Weight Decay* [Krogh and Hertz, 1991]. It is a well understood method from classical regression models, where estimation variance is reduced by regularising large weights or regression coefficients. Despite its easy interpretation for these classical models, its application to neural networks is not completely understood.

### A.2  Dropout

*Dropout* [Srivastava et al., 2014] is another commonly used regularisation method that works by overlaying noise during the learning procedure. The output from each neuron in the network is discarded at random during training according to a set distribution. Intuitively, by adding this noise, the network is forced to not rely on any given path through the network which could be used to memorise data.

### A.3  Learning rate algorithms and optimisers

There exist many optimisers other than standard SGD that also offer regularisation benefits. Optimisers such as Adam [Kingma and Ba, 2015] set an adaptive learning rate by using estimations of first and second order terms. On the other simple but effective learning rate algorithms such as *Cyclical Learning Rate* [Li and Yang, 2020] simply adjust the learning rate in a more nuanced way that leads to a better learning procedure without using any past training information.

### A.4  Loss gradient norm penalties

Loss gradient norm penalties add an additional penalty term to the loss: the norm of the gradient of the loss itself. This will lead the gradients of the gradients being calculated in the final back-propagation step, hence why it is commonly referred to as double back-propagation. First, a back-propagation pass must be done to obtain the gradients, then a second pass to evaluate the weight updates with the gradient norm included. Two versions exist: one which calculates the gradient of the loss with respect to the input data [Drucker and Le Cun, 1991, Hoffman et al., 2019] and another that instead does so with respect to the parameters [Barrett and Dherin, 2021, Zhao et al., 2022]:

$$\tilde{\mathcal{L}} = \mathcal{L} + ||\nabla_x \mathcal{L}||$$
$$\tilde{\mathcal{L}} = \mathcal{L} + ||\nabla_\theta \mathcal{L}||.$$

We note one important detail in the relation of the latter, the loss-parameter gradient norm, to the MGS trace penalty:

$$||\nabla_\theta \mathcal{L}|| = ||\nabla_\theta f \cdot \nabla_f \mathcal{L}|| \leq ||\nabla_\theta f|| \, ||\nabla_f \mathcal{L}||.$$

Therefore, $||\nabla_\theta L||$ is proportional to the MGS trace penalty as $||\nabla_\theta f||_2^2 = \sqrt{\operatorname{tr} K_\theta}$. However depending on how the second term $||\nabla_f \mathcal{L}||$ behaves the two might not be coupled at all. This might explain why in some instances in the experiments, the loss-parameter gradient penalty performed significantly worse. However, it is still the closest to directly optimising MGS and when the loss-model gradient is well behaved, can perform just as well as MGS. Still, the issue arises as to how this penalty performs when any minima is approached, as the second term will diminish, causing the whole penalty to lose influence over training. This might then lead to the model starting to overfit after landing in a minima, which we have observed in our experiments.

### A.5 Functional regularisation

Recently it has also been shown that *Kernel Ridge Regression* can be applied to neural networks. This places a penalty on function complexity by trying to minimise the norm of the function given by its *Reproducing Kernel Hilbert Space* (RKHS) (e.g. [Bietti et al., 2019]). For neural networks, the NTK can be used to define such a RKHS, however the actual function realised by the network is not necessarily close to its RKHS counterpart. Additionally, in our experience, these penalties are the most costly to calculate and were not efficient enough even to run on small networks.

## B Connection to Neural Tangent Theory

It should most certainly be noted that the MGS kernel is indeed the same as the *empirical Neural Tangent Kernel* $\hat{\Theta}$ which is derived from the *Neural Tangent Kernel* (NTK) itself. From the NTK literature, the NTK describes the evolution of the network function $f_\theta$ in function space in a kernel gradient descent setting. Therefore we can also make use of findings that have been made from the NTK perspective.

A observation made by about the empirical NTK is the notion of "feature learning speed" [Jacot et al., 2018, Baratin et al., 2021, Ronen et al., 2019]. Neural networks seem to learn functions/features of increasing complexity starting with low frequency content first and then sequentially higher frequency information as training progresses. The spectral decomposition of the empirical NTK describes the primary directions of learning with large eigenvalues corresponding to learning directions with low complexity which is also where convergence is the fastest. While most of the ideas from the NTK side are formulated by viewing the network as learning in function space, it can be useful to see what practical implications this has on the model gradients and therefore MGS.

Due to the NTK growing in popularity there exist a number of performant libraries for calculating the empirical NTK. We make use of *Neural Tangents* [Novak et al., 2019] and *Fast Finite Width NTK* [Novak et al., 2022] libraries which are built on *JAX* [Bradbury et al. [2018]].

## C Using MGS metrics for adversarial sample detection

The authors Martin and Elster [2020] used $\text{tr}\, K_\theta$ for detection of adversarial samples when running a neural network. They do this by calculating the trace of what they call the *approximate Fisher Information Matrix* (FIM) which is the actually same as $K_\theta$ in our case. They observed that if adversarial samples were present in the input, then the trace of the approximate FIM would be large. Their reasoning behind this behaviour, however, was deduced from the properties of the complete FIM. From the gradient similarity perspective it also makes sense. Adversarial samples are intended to look very similar to trained samples, but corrupt the model by causing it to produce different outputs. Therefore, if for adversarial data the trace is large, then the network considers that data to be very different from existing data. This also makes sense as that is one of the objectives of adversarial attacks in the first place.

## D Experiments

All experiments using a FCN were run locally on a laptop with a Nvidia Quadro T2000 graphics card. For convolutional networks, the laptop graphics card was not sufficient so instead a single Nvidia K80 was used on another machine.

### D.1 Experiment details

**Classification: MNIST**  A corrupted version (motion blur) of MNIST was used from the MNIST-C dataset [Mu and Gilmer, 2019]. Each method was tuned using a simple grid search over a set of parameters. The scenario in which the tuning took place was using 6000 randomly stratified sampled training points and 50% of the target labels were randomly flipped. Each method was trained for 100 epochs and also run 5 times over, with differently randomly sample data each time. The final test loss was used to determine the optimal parameter. The architectures tested included a standard fully connected architecture with 6 hidden layers, 300 neurons wide, and ReLU activations

as well as a LeNet-5 style convolutional network. Cross-entropy with softmax was used as the loss function. In the case of Dropout, a Dropout layer was added between all fully connected layers in both architectures as well as a Batchnorm layer between convolutional layers.

After tuning, each method retained the same parameters and was then run over a large testbench, for 400 epochs for both networks, where the training size and label noise were varied between 250 and 9000 training samples and 0% to 100% respectively. The test accuracy was calculated on the predetermined test sample found in MNIST and not the data that was left over from the training sample. Each scenario was run 10 times, again with new network initialisations and sampled training data. A slow exponential decay was used for the learning rate starting at 0.1 and a batch size of 32 was used.

Each of the metrics, including test accuracy, were sampled at 100 evenly spaced points during training, so as to not run too slowly. The final test accuracy was calculated as an average of the previous 5 metric samples and then averaged between the training runs. The max value was calculated as the max of the averaged runs.

**Regression: Facial Keypoints**    A similarly corrupted version (motion blur) of the Facial Keypoints dataset [1] was used. The same architecture types were used as in the MNIST experiment. Standard MSE loss is used. Unlike the MNIST experiment, a grid of scenarios was not run for this problem. Instead the training size was held fixed at 30% of the total data, while the amount of noise was changed. Gaussian noise was added to the target coordinates with a scale indicated by the "noise" column in the tables. Each method was tuned with a noise scale of 15 for 50 epochs.

Each scenario was then run for 400 epochs in the case of the FCN architecture and 50 epochs for the LeNet-5 architecture. This was due to the larger size of the input and output, causing the convolutional network to run much slower. A batch size of 128 was used for both architectures in an attempt to lessen the training wall-clock time. This could explain why MGS in some cases achieved slightly better performance for the FCN architecture over the LeNet-5 architecture as it was allowed to train for longer.

The metrics were tracked in the same way as the MNIST experiment, however test MSE was used instead of test accuracy.

**Training parameter robustness**    To test robustness with regards to different training parameters, a benchmark was run where each method is first tuned using the same scenario. Then, the tuned scenario is used as a starting point from which each of the 5 training parameters were changed independently. Both larger and smaller parameter values were selected to test as broad of a spectrum as possible. For each parameter changed, multiple runs were performed, again by resampling the data and using a different initialisation. The final performance results were aggregated, with the quantiles being drawn, based on all the results. The test was run using the same classification problem based on the corrupted MNIST dataset described previously. Specifically, each method was tuned using 3000 training samples, 50% label noise, and for 100 epochs with all the other settings such as learning rate the same as per section D.1.

## D.2    Results for MNIST and Facial Keypoints for a FCN network

---

[1]https://www.kaggle.com/competitions/facial-keypoints-detection/

Table A.1: **Test accuracy for corrupted MNIST dataset using a FCN architecture.** Noise column represents percentage of training labels that have been randomly flipped. MGS is able to handle large amounts of noise. If the max accuracy is compared to the final accuracy, MGS is also the most consistent and is not susceptible to overfitting. The FCC results are worse than those of LeNet-5.

| Noise | Unregularised | Dropout | Weight | Loss grad. | MGS |
|---|---|---|---|---|---|
| 0% | 82.1 ±1.9 (82.3) | 79.5 ±1.2 (80.0) | 72.8 ±9.7 (79.6) | 74.9 ±22.0 (79.8) | **84.6** ±1.2 (84.8) |
| 30% | 59.5 ±3.7 (75.3) | 73.7 ±2.5 (75.5) | 58.7 ±4.0 (74.6) | 70.3 ±4.3 (77.4) | **77.7** ±1.8 (79.4) |
| 60% | 36.9 ±2.8 (65.6) | 55.6 ±4.8 (57.1) | 10.0 ±0.5 (63.1) | 41.2 ±6.4 (43.4) | **67.8** ±3.4 (70.0) |
| 80% | 24.8 ±2.7 (53.0) | 26.5 ±3.7 (29.6) | 10.1 ±0.6 (41.3) | 12.1 ±2.5 (14.8) | **51.4** ±3.8 (55.0) |

Table A.2: **Test loss for the corrupted Facial Keypoints dataset using a FCN architecture.** MGS is the only regulariser that is able to attain convergence with good performance.

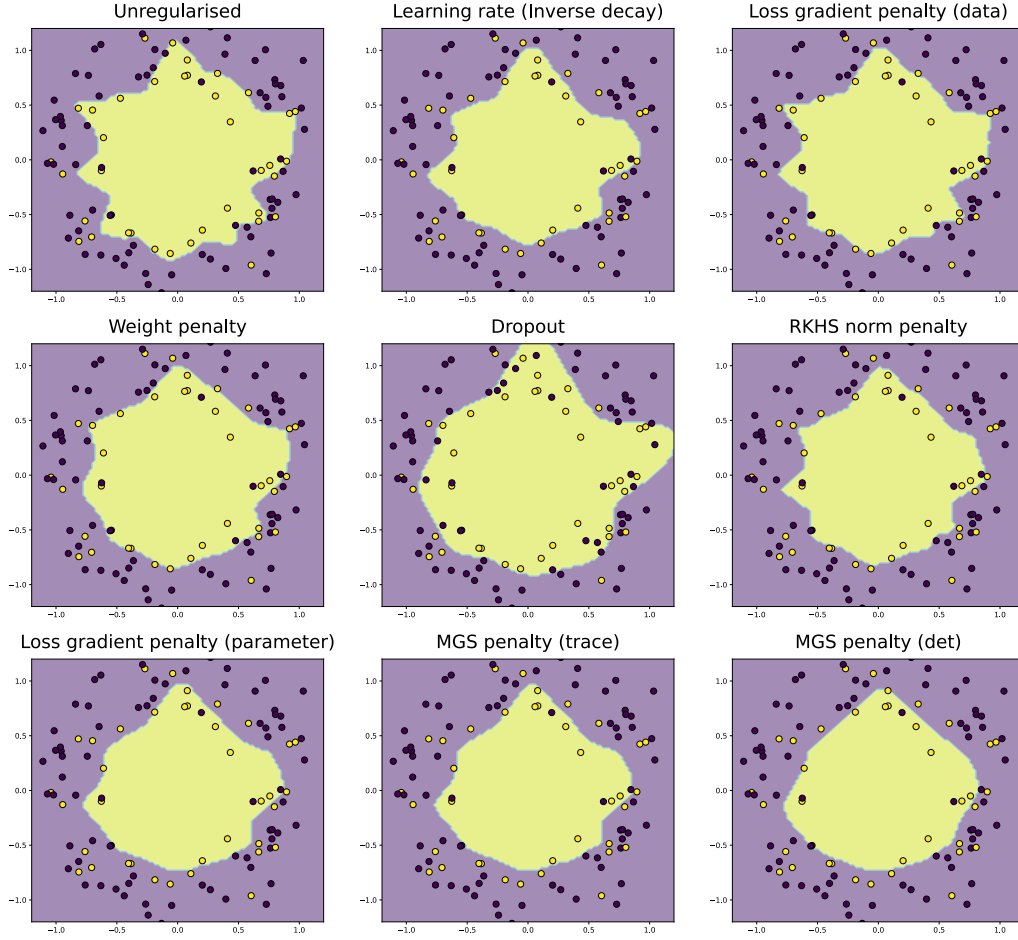| Noise | Unregularised | Dropout | Weight | Loss grad. | MGS |
|---|---|---|---|---|---|
| 0 | 281.1 ±264.6 (181.5) | 208.2 ±220.9 (74.3) | 276.4 ±262.3 (180.2) | 337.9 ±303.8 (184.7) | **10.6** ±0.3 (10.5) |
| 10 | 305.7 ±344.7 (176.8) | 240.2 ±291.1 (191.7) | 305.2 ±342.2 (175.0) | 335.3 ±332.1 (208.2) | **11.0** ±0.5 (10.7) |
| 20 | 325.2 ±371.7 (158.6) | 287.1 ±321.9 (88.0) | 328.5 ±379.4 (157.7) | 353.1 ±382.7 (174.2) | **13.6** ±1.3 (12.8) |
| 30 | 390.8 ±453.2 (162.6) | 210.4 ±270.9 (91.7) | 393.0 ±454.2 (162.2) | 411.8 ±442.8 (175.6) | **22.4** ±2.7 (20.9) |

# E    Two circles problem decision boundaries

Figure A.1: Decision boundaries for all regularisation methods.
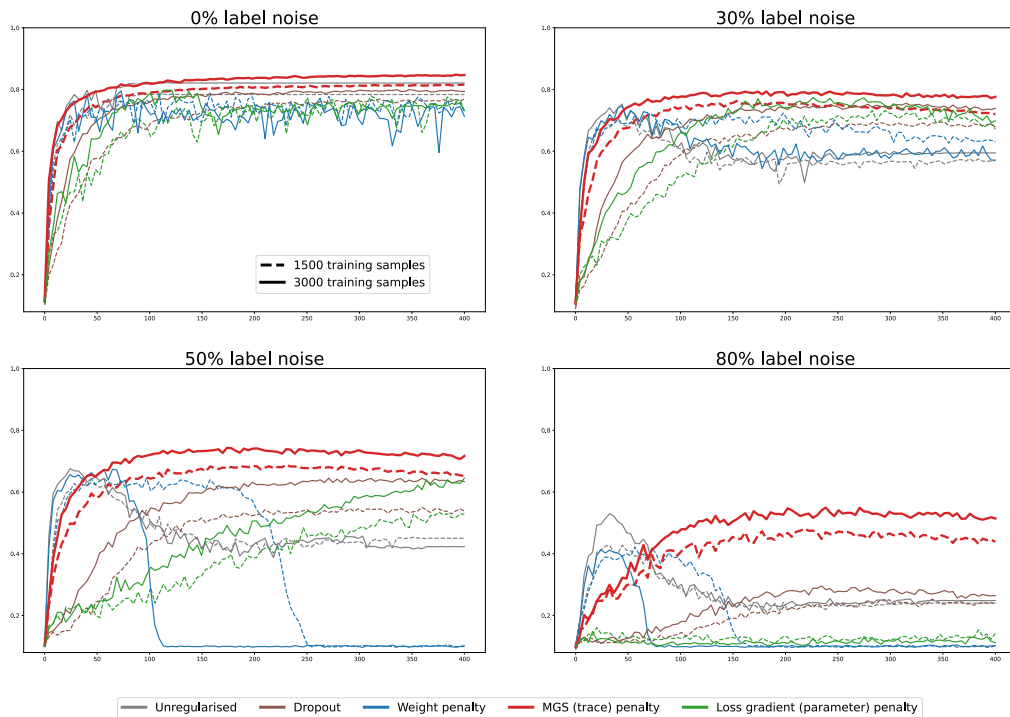
# F MNIST testbench results

Figure A.2: **Test accuracy during training using a FCN architecture.** Test accuracy plotted during training for two training sizes and different levels of label noise.
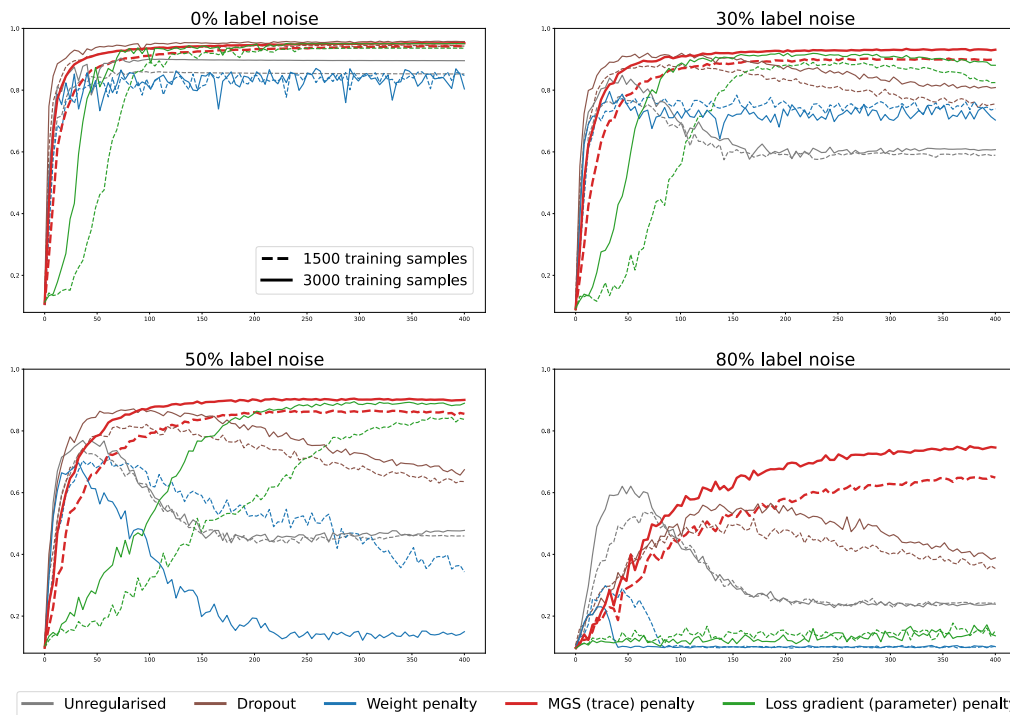


Figure A.3: **Test accuracy during training using a LeNet-5 architecture.** Test accuracy plotted during training for two training sizes and different levels of label noise.
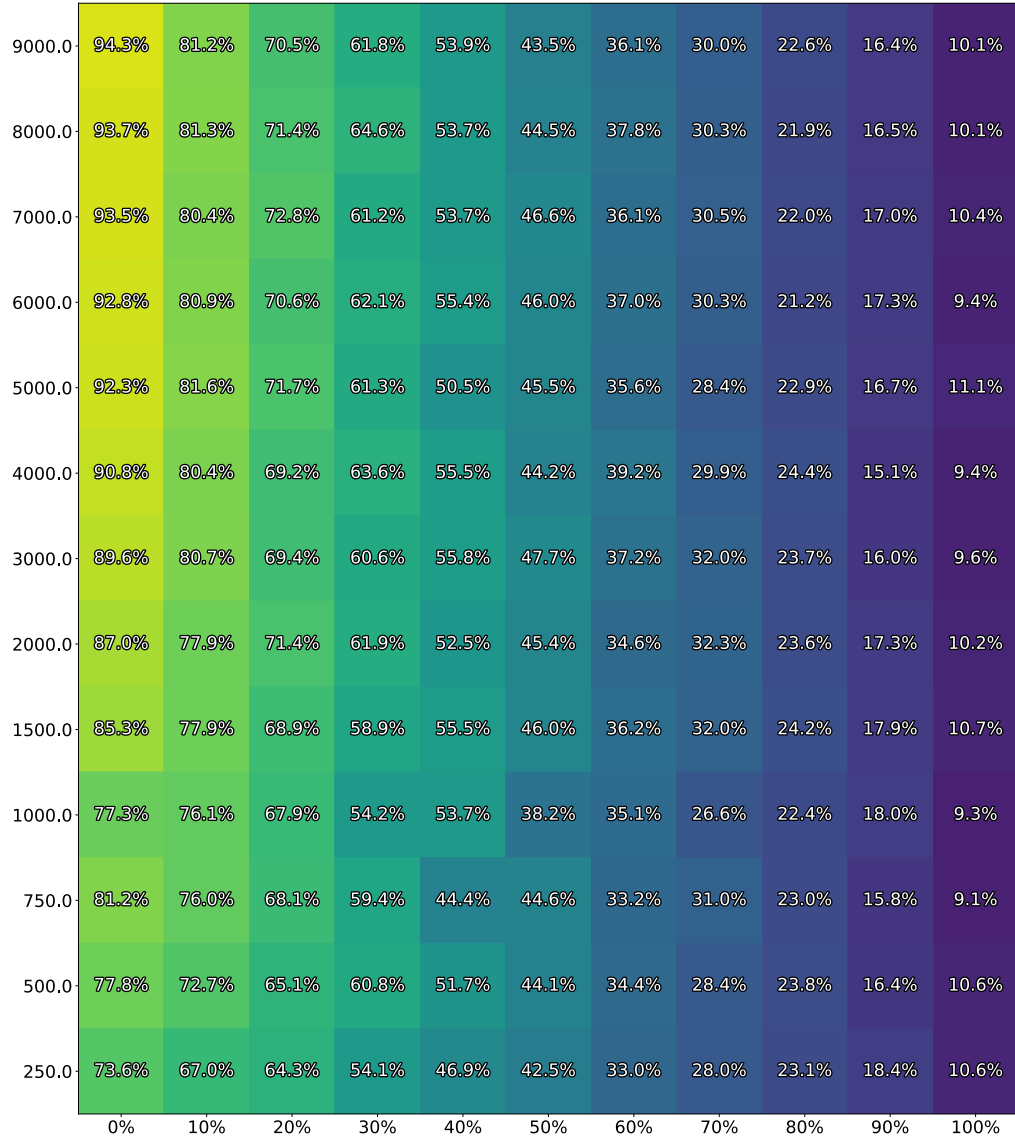
Figure A.4: **Unregularised LeNet-5**. Average final test accuracy for MNIST. Label noise on the x-axis and training size on the y-axis.

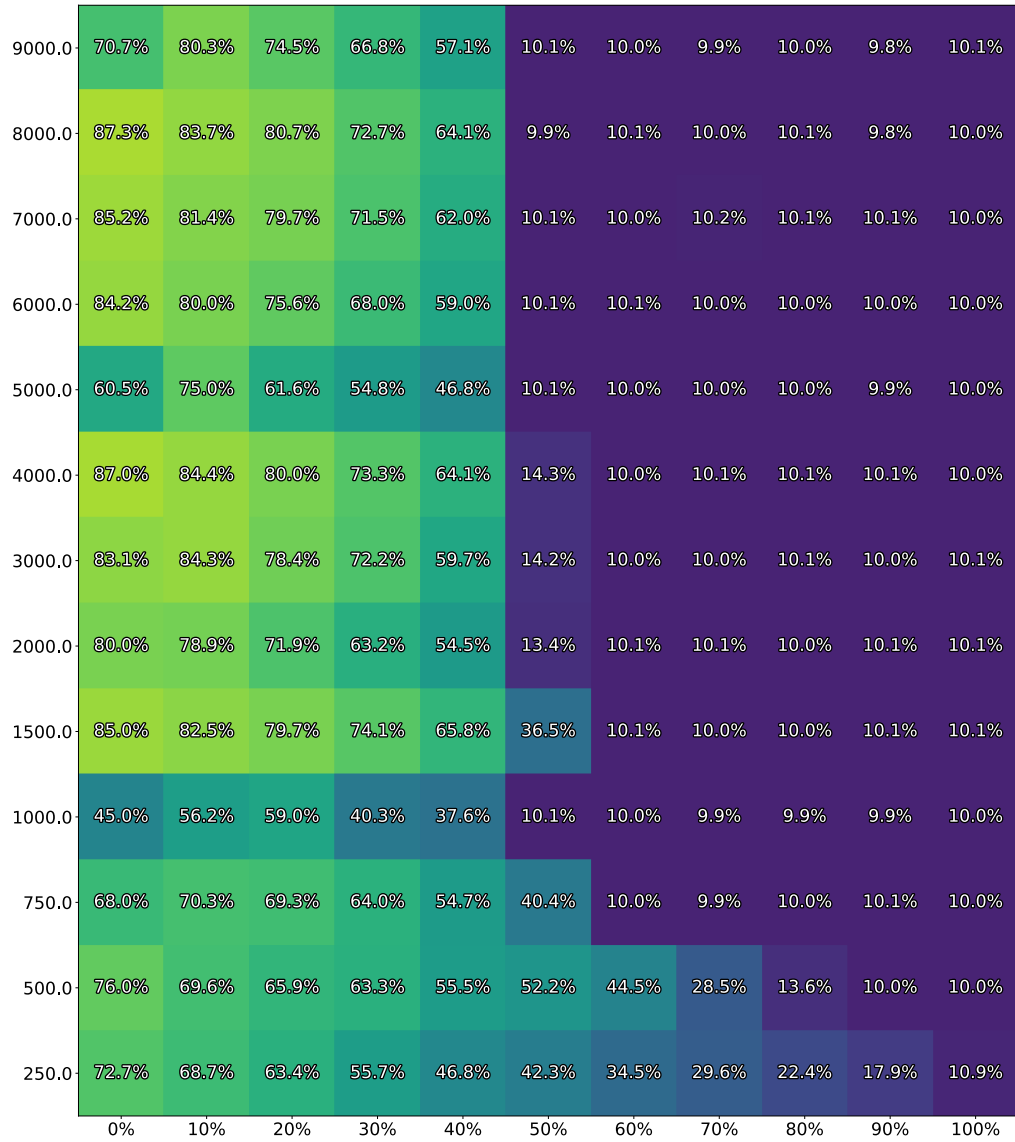| Training size | 0% | 10% | 20% | 30% | 40% | 50% | 60% | 70% | 80% | 90% | 100% |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 9000.0 | 70.7% | 80.3% | 74.5% | 66.8% | 57.1% | 10.1% | 10.0% | 9.9% | 10.0% | 9.8% | 10.1% |
| 8000.0 | 87.3% | 83.7% | 80.7% | 72.7% | 64.1% | 9.9% | 10.1% | 10.0% | 10.1% | 9.8% | 10.0% |
| 7000.0 | 85.2% | 81.4% | 79.7% | 71.5% | 62.0% | 10.1% | 10.0% | 10.2% | 10.1% | 10.1% | 10.0% |
| 6000.0 | 84.2% | 80.0% | 75.6% | 68.0% | 59.0% | 10.1% | 10.1% | 10.0% | 10.0% | 10.0% | 10.0% |
| 5000.0 | 60.5% | 75.0% | 61.6% | 54.8% | 46.8% | 10.1% | 10.0% | 10.0% | 10.0% | 9.9% | 10.0% |
| 4000.0 | 87.0% | 84.4% | 80.0% | 73.3% | 64.1% | 14.3% | 10.0% | 10.1% | 10.1% | 10.1% | 10.0% |
| 3000.0 | 83.1% | 84.3% | 78.4% | 72.2% | 59.7% | 14.2% | 10.0% | 10.0% | 10.1% | 10.0% | 10.1% |
| 2000.0 | 80.0% | 78.9% | 71.9% | 63.2% | 54.5% | 13.4% | 10.1% | 10.1% | 10.0% | 10.1% | 10.1% |
| 1500.0 | 85.0% | 82.5% | 79.7% | 74.1% | 65.8% | 36.5% | 10.1% | 10.0% | 10.0% | 10.1% | 10.1% |
| 1000.0 | 45.0% | 56.2% | 59.0% | 40.3% | 37.6% | 10.1% | 10.0% | 9.9% | 9.9% | 9.9% | 10.0% |
| 750.0 | 68.0% | 70.3% | 69.3% | 64.0% | 54.7% | 40.4% | 10.0% | 9.9% | 10.0% | 10.1% | 10.0% |
| 500.0 | 76.0% | 69.6% | 65.9% | 63.3% | 55.5% | 52.2% | 44.5% | 28.5% | 13.6% | 10.0% | 10.0% |
| 250.0 | 72.7% | 68.7% | 63.4% | 55.7% | 46.8% | 42.3% | 34.5% | 29.6% | 22.4% | 17.9% | 10.9% |

Figure A.5: **LeNet-5 regularised with weight penalty**. Average final test accuracy with label noise on the x-axis and training size on the y-axis.
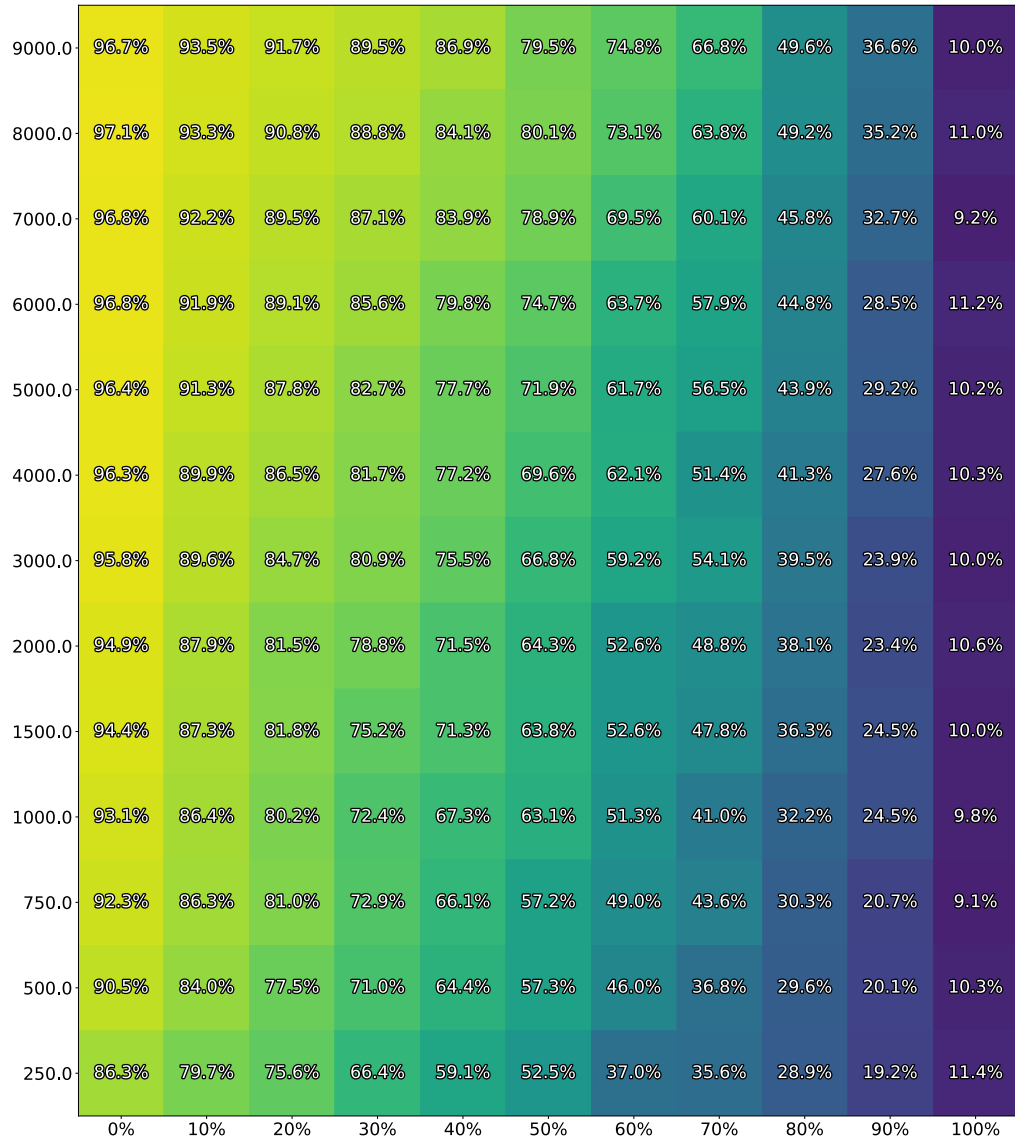
Figure A.6: **LeNet-5 regularised with dropout**. Average final test accuracy with label noise on the x-axis and training size on the y-axis.
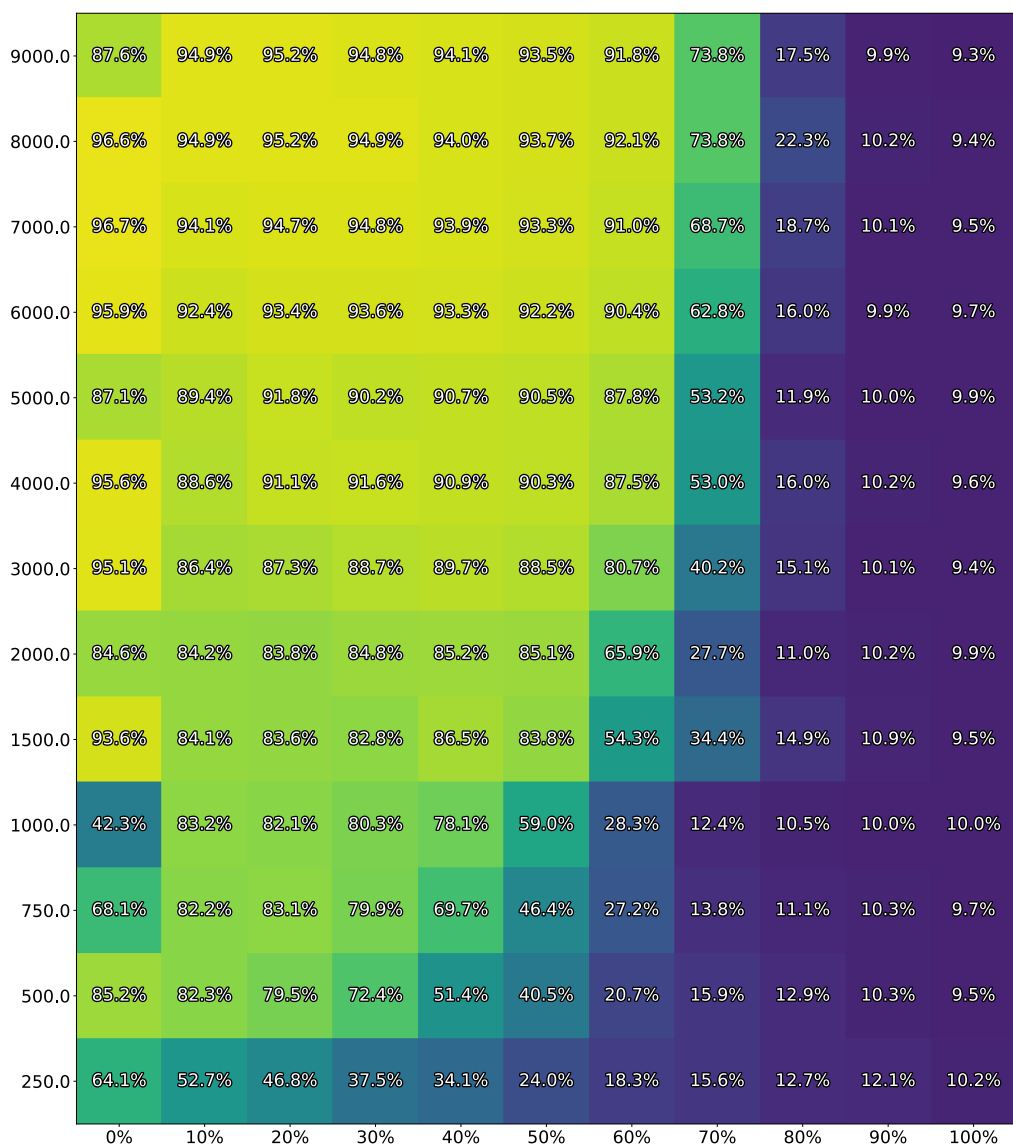
Figure A.7: **LeNet-5 regularised with loss-gradient penalty (parameter)**. Average final test accuracy with label noise on the x-axis and training size on the y-axis.
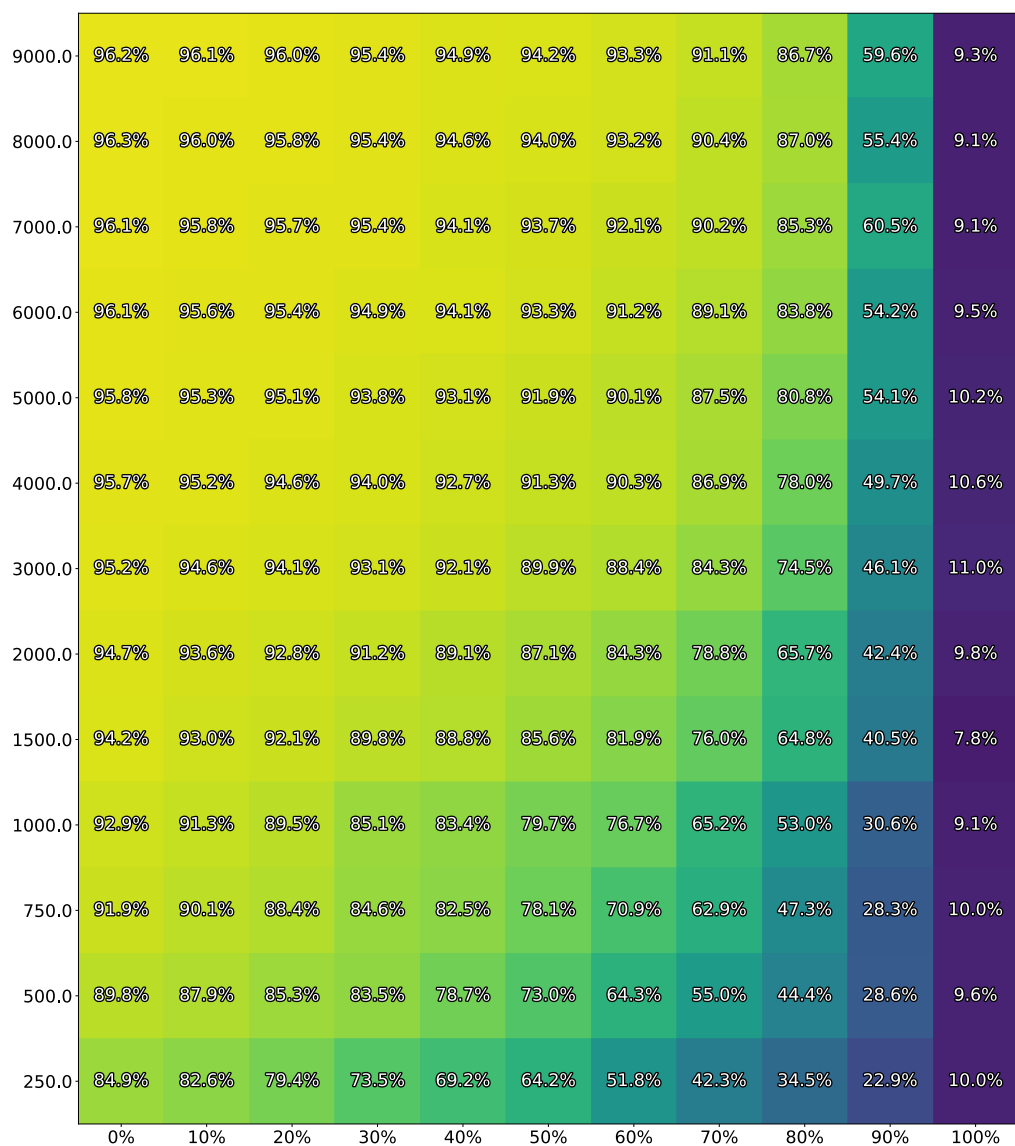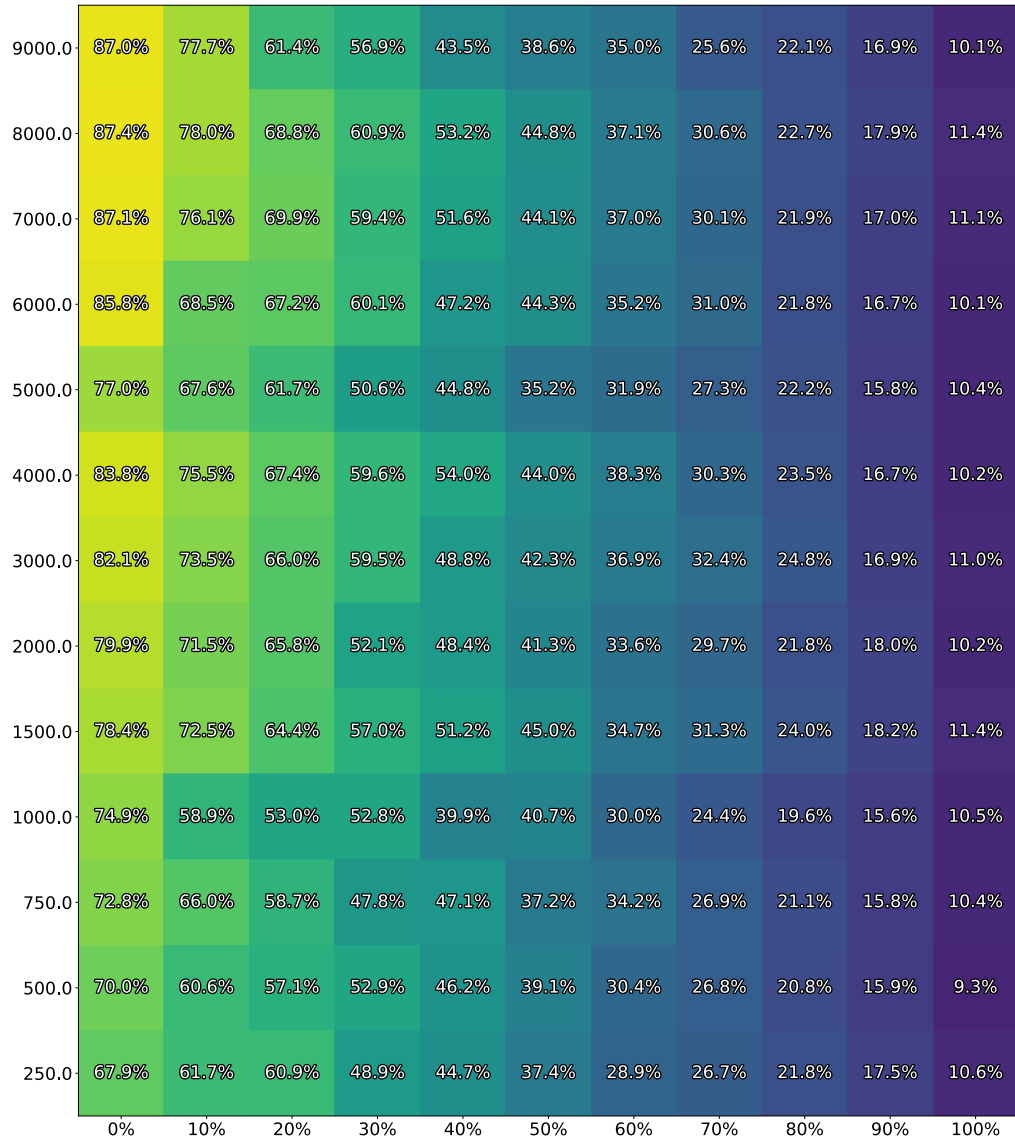
Figure A.8: **LeNet-5 regularised with MGS penalty (trace)**. Average final test accuracy with label noise on the x-axis and training size on the y-axis.

Figure A.9: **Unregularised FCN**. Average final test accuracy for MNIST. Label noise on the x-axis and training size on the y-axis.
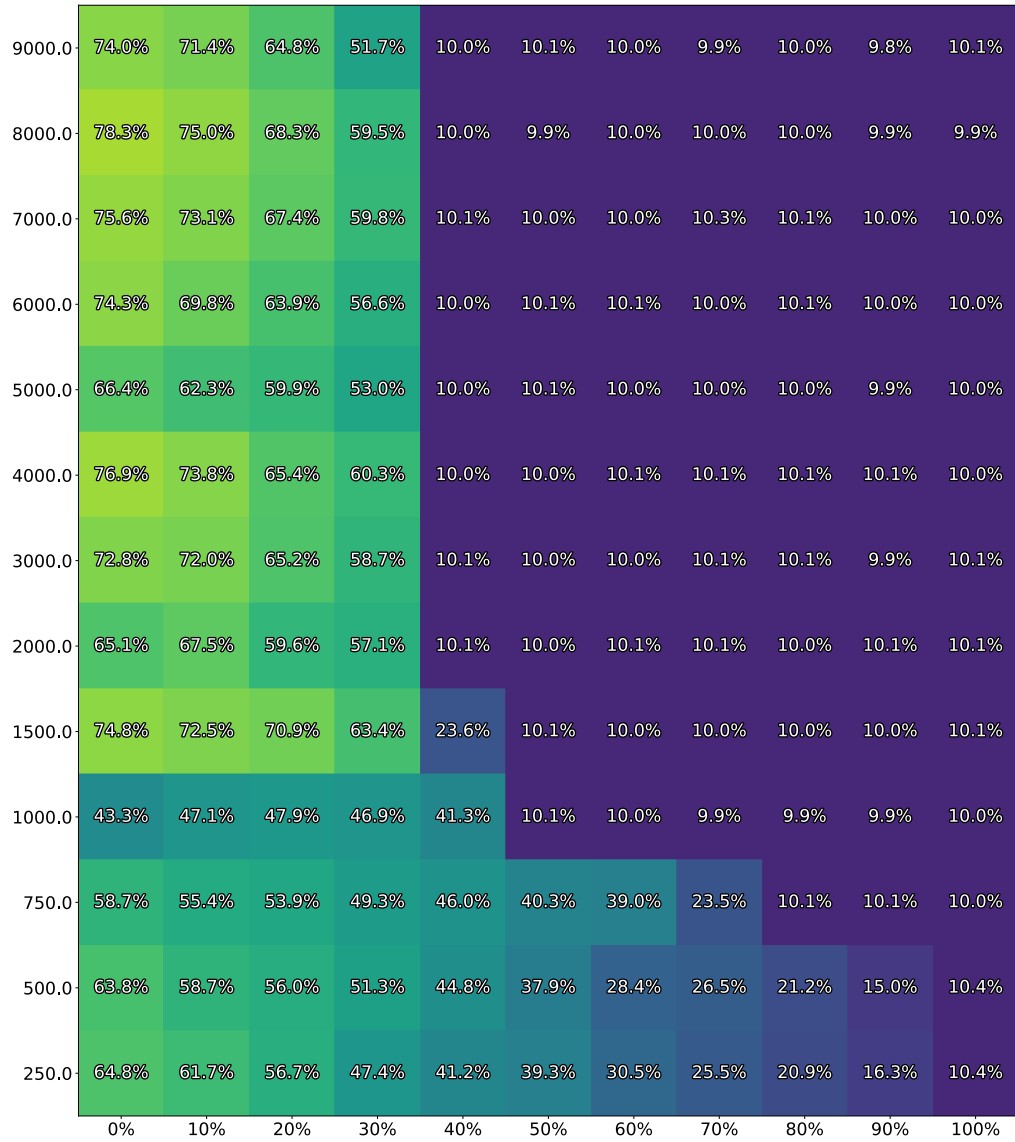
Figure A.10: **FCN regularised with weight penalty**. Average final test accuracy with label noise on the x-axis and training size on the y-axis.
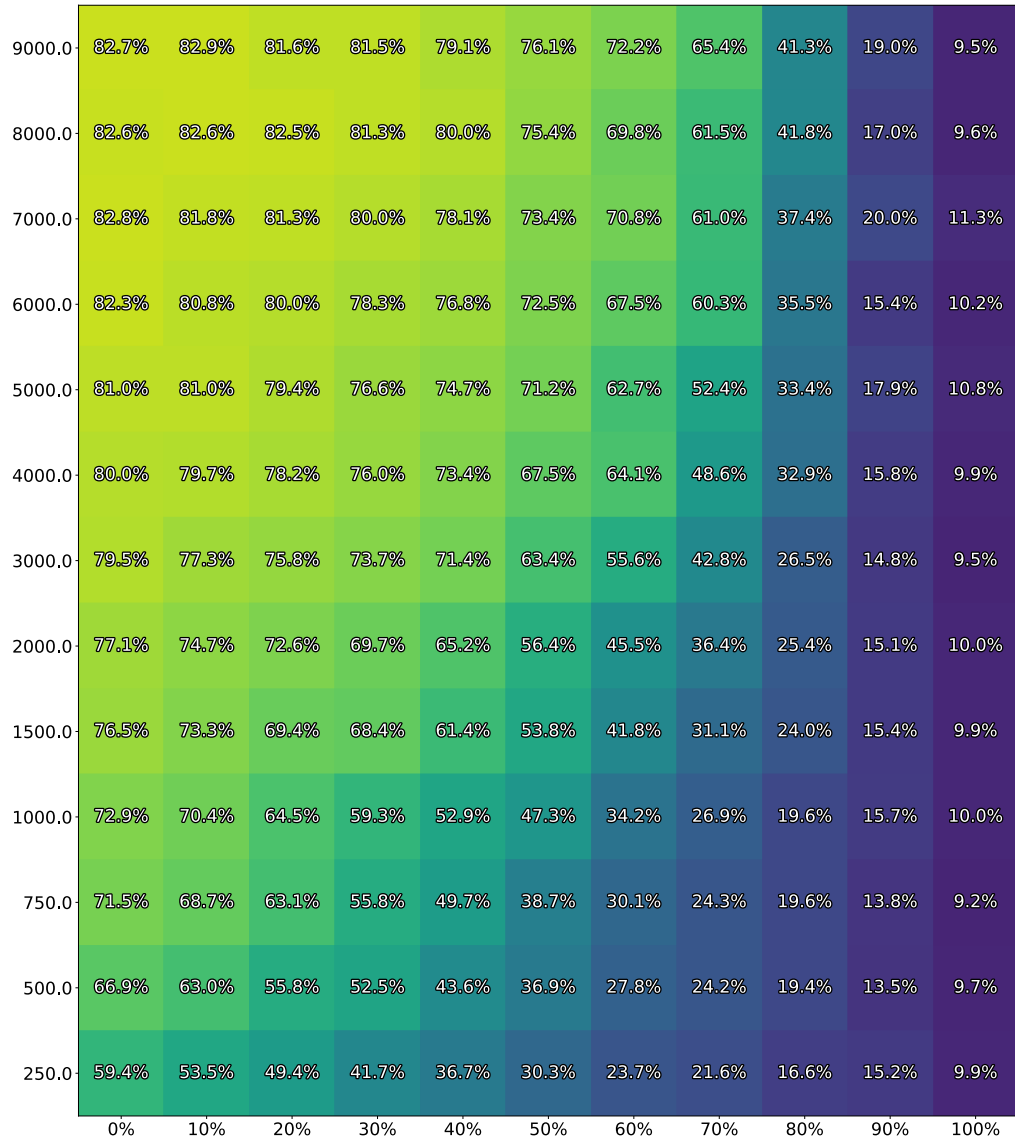
Figure A.11: **FCN regularised with dropout**. Average final test accuracy with label noise on the x-axis and training size on the y-axis.
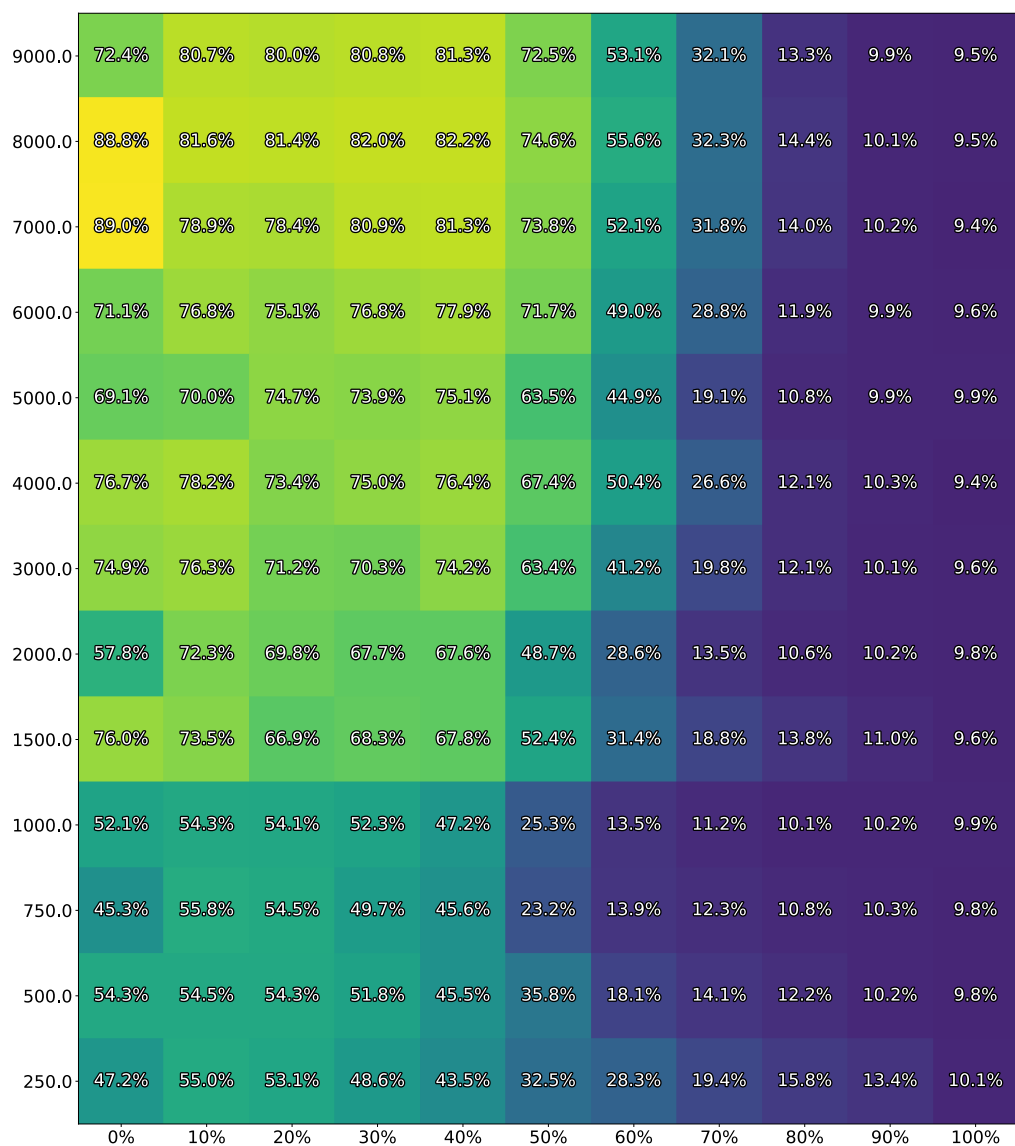
Figure A.12: **FCN regularised with loss-gradient penalty (parameter)**. Average final test accuracy with label noise on the x-axis and training size on the y-axis.
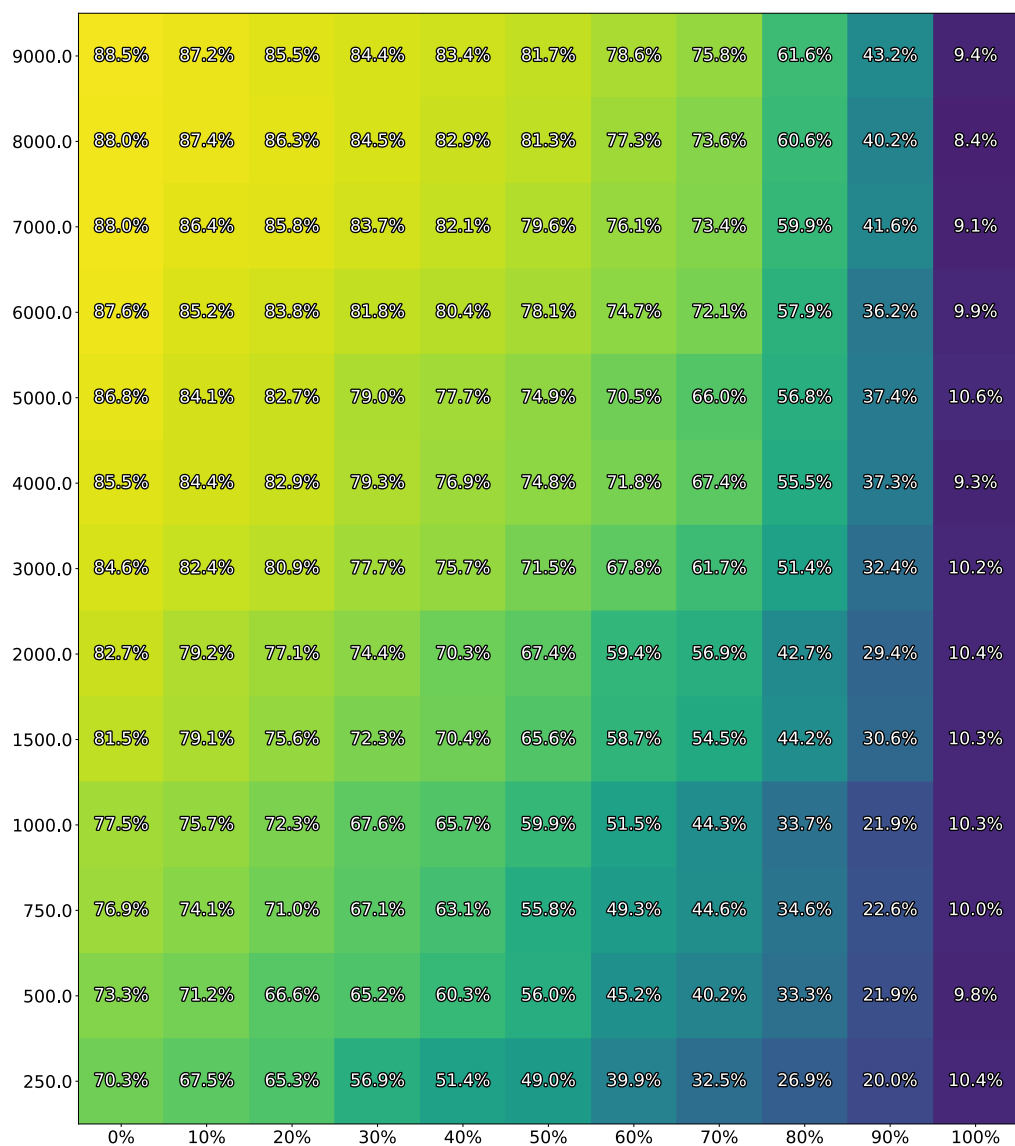
Figure A.13: **FCN regularised with MGS penalty (trace)**. Average final test accuracy with label noise on the x-axis and training size on the y-axis.

# G  Calculation of $K_\theta$

## G.1  Vector output models

In section 2 $K_\theta$ is defined for a scalar output model and therefore is a matrix of size $n \times n$. In the case of a vector output model $K_\theta$ becomes a tensor of size $n \times n \times q \times q$ where $q$ is the number of outputs. The reasoning about MGS still holds. However, the calculations become more involved as one not only has to consider the model update for a single output (given by the diagonal matrices $K_\theta^{i,j}, i = j$) but also how an update for one output affects other outputs (given by the off diagonal matrices $K_\theta^{i,j}, i \neq j$). Therefore, it is easier to instead combine all the outputs by concatenating their model gradients and then calculate a single kernel $K_\theta$ of size $np \times np$, as is common in the NTK literature [Lee et al., 2019]:

$$\nabla_\theta f = [\nabla_\theta f_1; \cdots ; \nabla_\theta f_q].$$

## G.2  Handling of mini-batches and large data sets

When using mini-batches, $K_\theta$ is no longer square and instead will be of size $n \times m$, where $m$ is the number of samples in the current batch.

Both this matrix and the $n \times n$ matrix from standard gradient descent can be infeasible to calculate for large data sets. Instead, an approximation is calculated using only the data from the current batch which will be of size $m \times m$. Throughout the entire article, $K_\theta$ is always calculated for the current batch and not the full dataset.

# H  Relation between spectrum of MGS kernel and gradient sample covariance matrix

**Proposition 1**  *The matrices $X^T X$ (Gram matrix) and $X X^T$ (scatter matrix) share the same non-zero eigenvalues.*

*Proof: Suppose that $\lambda$ is a non-zero eigenvalue of $X^T X$ with the associated eigenvector u.*

*Then:*

$$X^T X u = \lambda u$$
$$X X^T X u = X \lambda u$$
$$X X^T (X u) = \lambda (X u)$$
$$X X^T \tilde{u} = \lambda \tilde{u}.$$

*Thus, $\lambda$ is an eigenvalue of $X X^T$, with the associated eigenvector $\tilde{u} = XU$.∎*

**Theorem 1**  *Separation theorem [Takane and Shibayama, 1991].*

*Let $M$ be a d-by-n matrix. Let two orthogonal projection matrices be $P_{left}$ and $P_{right}$ of size d-by-d and n-by-n respectively.*

*Then:*
$$\sigma_{j+t}(M) \leq \sigma_j(P_{left} M P_{right}) \leq \sigma_j(M),$$
*where $\sigma_j(\cdot)$ denotes the j-th largest singular value of the matrix, while $t = d - r(P_{left}) + n - r(P_{right})$ and $r(\cdot)$ is the rank of the matrix.*

**Theorem 2**  *Let $C$ and $\bar{C}$ be the scatter matrix and its centred counterpart, then their eigenvalues are interlaced, such that [Honeine, 2014]:*

$$\lambda_{j+1} \leq \bar{\lambda}_j \leq \lambda_j.$$

*Proof: Let $M = X$, $P_{left}$ being the d-by-d identity matrix and $P_{right} = (I - \frac{1}{n} 11^T)$ which is the n-by-n centering matrix.*

*With $r(P_{left}) = d$ and $r(P_{right}) = n - 1$, it follows from the Separation theorem 1 that:*

$$\sigma_{j+1}(X) \leq \sigma_j(X - (I - \frac{1}{n}11^T)) \leq \sigma_j(X)$$

*Furthermore it is well known that the eigenvalues of $C$ are equal to the square roots of the singular values of $X$.*∎

By using proposition 1 and theorem 2 we see that the eigenvalues of the MGS kernel, which is the Gram matrix for the model-gradients, and the centred scatter matrix (sample-covariance matrix) will be interlaced.