
Easy incremental learning methods to consider for commercial fine-tuning applications

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Fine-tuning deep learning models for commercial use cases is growing exponentially as more and more companies are adopting AI to enhance their core products
2 and services, as well as automate their diurnal processes and activities. However,
3 not many countries like the U.S. and those in Europe follow quality data collection
4 methods for AI vision or NLP related automation applications. Thus, on many
5 of these kinds of data, existing state-of-the-art pre-trained deep learning models
6 fail to perform accurately, and when fine-tuning is done on these models, issues
7 like catastrophic forgetting or being less specific in predictions as expected occur.
8 Hence, in this paper, simplified incremental learning methods are introduced to be
9 considered in existing fine-tuning infrastructures of pre-trained models (such as
10 those available in huggingface.com) to help mitigate the aforementioned issues
11 for commercial applications. The methods introduced are: 1) Fisher Shut-off, 2)
12 Fractional Data Retention and 3) Border Control. Results show that when applying
13 these methods on vanilla pre-trained models, the models are in fact able to add more
14 to their knowledge without hurting much on what they had learned previously.
15

16 1 Introduction

17 Many companies and organizations today are adopting AI in automation, automating their daily
18 processes and activities, as well as offering them in their core products and services. Automation
19 has traditionally been in the industry for many years, as a means for which economics of scale could
20 be achieved so as to remain competitive in the market. Now with AI, more and more intelligence is
21 being brought into automation, and in countries like India, organizations are beginning to adopt AI
22 for this particular purpose.

23 With recent advancements in AI vision and NLP models such as the GPT-3, Jurassic-1, and so on,
24 organizations today are using AI for 1) Document Reading and Understanding, 2) Online Proctoring,
25 3) Chatbots, 4) Intelligent Information Parsing and other application related process automations.
26 Given these use cases, AI solutions need to be specific to their processes, but yet be an addition to
27 their generally known formats. This in a sense, is more like making use of a human employee who
28 has some kind of general education on various tasks or processes but still is required to learn the
29 companies counterparts well and in detail before he/she is allowed to execute them. These processes
30 can include between, reading customer emails for entering relevant information about their product
31 requirements onto a structured database, to understanding various types of printed documents for
32 information parsing, and to identifying newer objects for either document filtering or malicious
33 activity detection.

34 For natural language related tasks, powerful models like the GPT-3 are now being widely used, but
35 they require good prompt engineering skills to get the best out of them. Also, given that they are
36 probabilistic models, the generated outputs can sometimes falter away from what is expected, and

37 this can become a problem when selling it to customers, because even the slightest faltering may not
38 be acceptable to them at all. Hence, to reduce this, more and more examples have to be provided in
39 the prompt, and this can come at a high cost not suitable for low cost of living countries like India.
40 The other workaround is to fine-tune the model on the new datasets, but this has epoch limitations
41 on how deeply it can fit on the new dataset without hurting the body of general knowledge it gained
42 earlier. Also, fine-tuning models like the GPT-3 comes at a very high cost now-a-days, and is no
43 more an option. This leaves the automation builders to use `huggingface.com` transformers instead.

44 In vision, although state-of-the-art pre-trained deep learning models are able to achieve human level
45 performance on a variety of inputs, they can only perform so in upto close to high quality inputs. If
46 the quality goes lower, they fail terribly. Not all organizations have a good quality data collection
47 process involved for applying automation, and this is ubiquitously the case in many parts of the world.
48 So it becomes quite difficult to sell AI as a human-level performer, and at this point AI becomes of
49 lesser use than it could potentially be.

50 Another approach typically used to resolve such problems is to employ transfer learning, which
51 typically involves replacing the last layers of the model with a new model to get the specific outputs
52 required. Some examples done in research are Too, et al. (2019), Dif & Elberrichi (2020), Alshalali
53 & Joysula (2018), Jung, et al. (2015), Qian, et al. (2021) and Vrbančič & Podgorelec (2020). While
54 this may not seem to be a problem with vision based tasks, it is definitely a problem with natural
55 language based tasks. This is because the final layers of the natural language models have all the vital
56 information of language structure that help with the language generative process. When this is to
57 be changed, catastrophic forgetting can happen. Catastrophic forgetting is a phenomenon in which
58 previously learned knowledge is lost partly by the application of new data for training. Also, with
59 vision based tasks, when the requirement is to just improve the performance on lower quality data,
60 transfer learning may not be the appropriate approach. Fine-tuning for these must involve the final
61 layers of the model which could inevitably lead to catastrophic forgetting on the higher quality inputs.

62 This brings the only solution towards **incremental learning**. This type of learning is all about
63 learning on newer datasets without having the side-effects catastrophic forgetting, and there has been
64 substantial amount of research done in this area. Luo, et al. (2020) summarizes all the work that has
65 happened in this area so far. There are several approaches to implementing incremental learning on
66 pre-trained models, some of which will be discussed in the forthcoming sections. In this paper, a
67 few of these approaches will be simplified for commercial applications along with novel intuitive
68 additions to further help the learning process. The paper introduces: 1) Fisher Shut-off which is a
69 simplification of the work done by Kirkpatrick, et al. (2017), 2) Fractional Data Retention which
70 adopts ideas from Castro, et al. (2018), and 3) Border Control which is an extension to the idea
71 outlined by Ren, et al. (2018) on reweighting examples by employing a method similar to Adaboost.
72 The last one is the novel addition as it formulates a different approach to retaining salient examples
73 for incremental learning. It is based on the work by Ruping (2001) on incremental learning with
74 SVMs. But since SVMs are too complex in the context on neural networks, a similar but simplified
75 approach is proposed.

76 The purpose of this work is to initiate the development of a new infrastructure for commercial
77 fine-tuning of pre-trained models with simplified incremental learning methods.

78 The rest of this paper proceeds as follows: Section 2 will provide a brief discussion on incremental
79 learning methods developed so far, followed by the proposal of simplified incremental learning
80 methods in Section 3. Section 4 will show sample results of the proposed methods on a vanilla
81 pre-trained model using a toy dataset. A toy dataset is used for the only purpose of providing
82 visualizations on the performance of the proposed methods. Nevertheless, these methods can be
83 extended on to real world datasets. The paper then concludes in Section 5 discussing steps forward
84 for implementation.

85 **2 Incremental Learning**

86 This section is a summary of the review published by Luo, et al. (2020). In this review, four different
87 types of strategies for incremental learning are highlighted, and every work published in this area
88 uses either one or more such strategies. Some examples are Castro, et al. (2018) and He, et al. (2020).
89 The four strategies are:

- 90 • Architectural
- 91 • Regularization
- 92 • Rehearsal
- 93 • Pseudo-Rehearsal

94 The following subsections will discuss these briefly.

95 2.1 Architectural Strategy

96 This strategy is similar to boosting techniques where multiple models are trained. But when used in
 97 the context of incremental learning, each model is trained on a different task separately. Then another
 98 meta-model that effectively selects which model to use for inference is trained. The work done by
 99 Poliker, et al. (2001) resembles this in many ways. In this work, multiple classifiers are trained with
 100 different training sets, and then a Adaboost style of ensemble learning is employed to combine the
 101 model outputs.

102 Another interesting work is by Rusu, et al. (2016) on Progressive Neural Networks (PNN). In this
 103 work, a neural network is trained sequentially on different tasks or training sets. However, each time,
 104 new neurons are added in each layer with new weights, and the weights of the previously learned
 105 neural network are frozen. Then, to prevent catastrophic forgetting, the outputs of each layer of the
 106 previous neural network on the earlier training set are used in addition to the new task or training
 107 set, when training the new layer neurons. The results on this type of incremental learning were quite
 108 encouraging that it set a new direction in the research of dynamically expanding networks that could
 109 make better use the neural networks capacity than the PNN. In fact, it will be seen later that the Fisher
 110 Shut-off method proposed in this paper inherently employs the idea of PNNs.

111 2.2 Regularization Strategy

112 In this strategy, as the name suggest, a regularization term is added in the loss function that measures
 113 the importance of old knowledge when learning on a new training set. The representative work done in
 114 this is Kirkpatrick, et al. (2017), whereby they introduce the concept of Elastic Weight Consolidation
 115 (EWC) by means of a Fisher Information Matrix. The EWC brings about the regularization term in
 116 the loss function as

$$R(w) = \sum_i \frac{\lambda}{2} F_i (w_i - w_{i,old})^2 \quad (1)$$

117 where F_i is the Fisher Information Matrix which suggests the importance of the i -th weight trained
 118 on the old (or previous) training set. Here, as one could speculate, the term Fisher Shut-off proposed
 119 in this paper actually derives itself from the Fisher Information Matrix, meaning that this matrix is
 120 used as the basis for shutting off the training of certain weights when training on a new set.

121 Another popular type of regularization strategy is Knowledge Distillation introduced by Hinton, et al.
 122 (2015). In this method, knowledge from an ensemble of models trained on different tasks (or training
 123 sets) separately are distilled into a smaller model that can be deployed much easily for inference.
 124 There are many huggingface.com transformers that are a product of such knowledge distillation.
 125 The distillation ensures that the smaller model holds all the knowledge of the ensemble, and that it
 126 can infer as good as it. Distillation is done by setting soft-targets on the smaller network from all the
 127 earlier training sets of the ensemble. The soft-targets are the output logits from the ensemble models
 128 on their respective trained datasets.

129 2.3 Rehearsal and Pseudo-Rehearsal Strategies

130 Rehearsal strategies in incremental learning make use of the earlier training sets when training a
 131 model on new tasks or training sets. This by far is the simplest of all incremental learning strategies
 132 that ensures catastrophic forgetting is prevented. The only issue is that when this strategy is used for
 133 deep learning models trained on large datasets, the training on new datasets could become extremely

134 slow and even time consuming before any fruitful results are achieved. Hence, newer research work in
 135 this area formulate methods for retaining only the most important data points to prevent catastrophic
 136 forgetting. The work done by Castro, et al. (2018) is an example of this. In this work, selection and
 137 removal mechanisms on data are introduced for assimilation into a memory network.

138 Talking about memory networks, the Pseudo-Rehearsal strategy involves training an additional data
 139 generator to generate the samples, the neural network was trained on earlier. Hence, newer research
 140 in this area involve GANs for data generation. Examples are Odena, et al. (2017) and Wu, et al.
 141 (2018).

142 3 Proposed Incremental Learning Methods

143 Commercial applications always require simplistic implementations of advanced methods no matter
 144 how complex they may be. Therefore, it is for this purpose alone this paper proposes some simplified
 145 methods for implementing incremental learning. As mentioned earlier in Section 1, these methods are:
 146 1) Fisher Shut-off, 2) Fractional Data Retention, and 3) Border Control. This section covers them in
 147 detail.

148 3.1 Fisher Shut-off

149 As mentioned in the previous section, the term Fisher Shut-off derives itself from the Fisher Infor-
 150 mation Matrix which weighs the importance of weights trained on previous datasets. Hence, in this
 151 sub-section, a brief overview of the details behind this matrix is covered with the help of Aich (2021).

152 Let \mathcal{D} represent a dataset coming from a stream of data for incremental learning. Then $p(w|\mathcal{D})$
 153 represents the model trained on data \mathcal{D} . This means that to train a model on a new dataset, the
 154 following posterior must satisfy:

$$p(w|\mathcal{D}_{new}) = \frac{p(\mathcal{D}_{new}|w)p(w|\mathcal{D}_{old})}{p(\mathcal{D}_{new})} \quad (2)$$

155 Note here that $p(w|\mathcal{D}_{old})$ is written in place of $p(w)$ because when \mathcal{D}_{new} is applied to the model, the
 156 weights w have already been trained with \mathcal{D}_{old} . Hence, given the model, $p(w|\mathcal{D}_{old})$, the log-likelihood
 157 loss on \mathcal{D}_{new} becomes,

$$\begin{aligned} \mathcal{L}_{\mathcal{D}_{new}}(w) &= \log(p(w|\mathcal{D}_{new})) \\ &= \log(p(\mathcal{D}_{new}|w)) + \log(p(w|\mathcal{D}_{old})) - \log(p(\mathcal{D}_{new})) \\ &\approx \log(p(\mathcal{D}_{new}|w)) + \log(p(w|\mathcal{D}_{old})) \end{aligned} \quad (3)$$

158 Here, the $\log(p(\mathcal{D}_{new}|w))$ equals the cross-entropy loss of the model on \mathcal{D}_{new} while $\log(p(w|\mathcal{D}_{old}))$
 159 is loss of the model on \mathcal{D}_{old} . To ensure that catastrophic forgetting does not occur on \mathcal{D}_{old} in its
 160 absence while training on \mathcal{D}_{new} , the loss on \mathcal{D}_{old} will have to be approximated using w alone. To do
 161 this, the Taylor's expansion on $\log(p(w|\mathcal{D}_{old}))$ is taken as,

$$\begin{aligned} \mathcal{L}_{\mathcal{D}_{old}}(w) &\approx \mathcal{L}(w)|_{\mathcal{D}_{old}} + \left(\frac{\partial \mathcal{L}(w)}{\partial w} \Big|_{\mathcal{D}_{old}} \right) + \frac{1}{2} (w - w|_{\mathcal{D}_{old}})^T \left(\frac{\partial^2 \mathcal{L}(w)}{\partial^2 w} \Big|_{\mathcal{D}_{old}} \right) (w - w|_{\mathcal{D}_{old}}) \\ &\approx \mathcal{L}(w)|_{\mathcal{D}_{old}} + \frac{1}{2} (w - w|_{\mathcal{D}_{old}})^T \left(\frac{\partial^2 \mathcal{L}(w)}{\partial^2 w} \Big|_{\mathcal{D}_{old}} \right) (w - w|_{\mathcal{D}_{old}}) \end{aligned} \quad (4)$$

162 since technically $\frac{\partial \mathcal{L}(w)}{\partial w} \Big|_{\mathcal{D}_{old}} = 0$, if the model is trained well on \mathcal{D}_{old} . Then, noting that the last term
 163 in (4) is equivalent to a regularization term, this term alone could be considered as the loss on \mathcal{D}_{old}

164 for preventing catastrophic forgetting. In doing so, the Fisher Information Matrix will equal to the
 165 Hessian, $\frac{\partial^2 \mathcal{L}(w)}{\partial^2 w} \Big|_{\mathcal{D}_{old}}$. This Hessian, \mathcal{H} , can be simply computed by the model gradients $\frac{\partial \mathcal{L}(w)}{\partial w} \Big|_{\mathcal{D}_{old}}$
 166 assuming that not all gradients are zero, as,

$$\mathcal{H} = \frac{\partial \mathcal{L}(w)}{\partial w} \Big|_{\mathcal{D}_{old}} \cdot \frac{\partial \mathcal{L}(w)}{\partial w} \Big|_{\mathcal{D}_{old}}^T \quad (5)$$

167 Doing so, and keeping only the diagonal terms, would imply that the model gradients are more
 168 than enough to weigh the important weights of the model trained on \mathcal{D}_{old} . Replacing (5) in (4) and
 169 substituting in (3) would give the loss on \mathcal{D}_{new} as,

$$\mathcal{L}_{\mathcal{D}_{new}}(w) \approx \log(p(\mathcal{D}_{new}|w)) + \frac{1}{2}(w-w|_{\mathcal{D}_{old}})^T \left(\frac{\partial \mathcal{L}(w)}{\partial w} \Big|_{\mathcal{D}_{old}} \cdot \frac{\partial \mathcal{L}(w)}{\partial w} \Big|_{\mathcal{D}_{old}}^T \right) (w-w|_{\mathcal{D}_{old}}) \quad (6)$$

170 which to an extent implies that if the model gradients on \mathcal{D}_{old} are absolutely zero, they get trained on
 171 \mathcal{D}_{new} without regularization, while those that are not, get regularized towards $w|_{\mathcal{D}_{old}}$.

172 This is what the proposed Fisher Shut-off exploits. In Fisher Shut-off, all weights of the model
 173 trained on \mathcal{D}_{old} that do **not** have absolute zero gradients get shut-off for training on \mathcal{D}_{new} , while the
 174 remaining that do take part. Also, since in practice *ReLU* functions are commonly used in deep
 175 learning models as the activation functions of the neurons, shutting off these weights becomes as
 176 simple as setting a condition. Figure 1 shows a sample performance of Fisher Shut-off on a regression
 177 model trained sequentially on mutually exclusive batches of data. These batches could represent the
 178 different tasks or training sets.

179 However, when it comes to classification, simple shut-off does not work completely. This is because,
 180 while in regression problems datasets could inherently employ some kind of piece-wise nonlinear fit
 181 in their distributions, the same cannot always be guaranteed in classification. Thus, in classification,

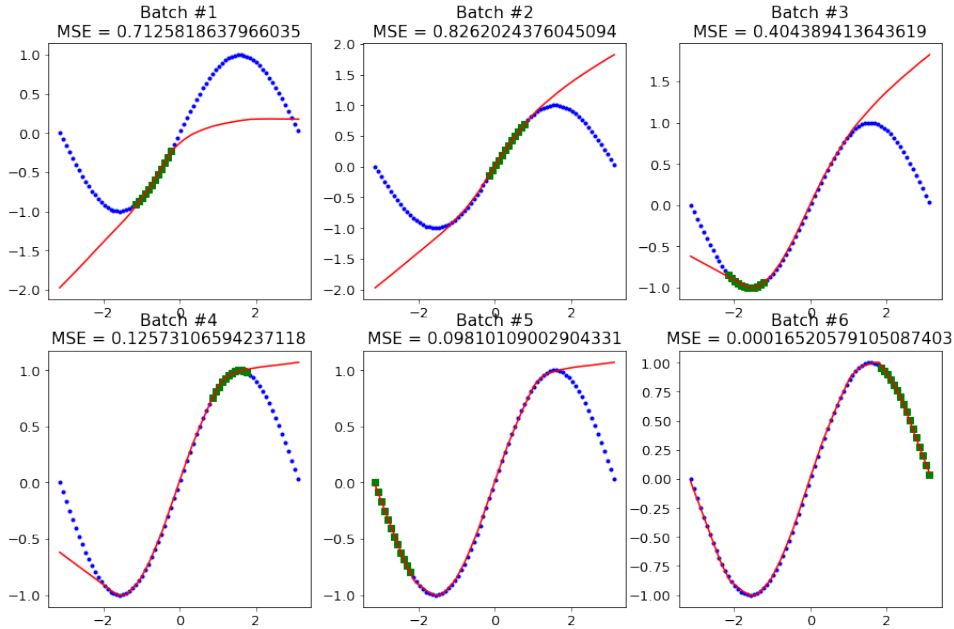


Figure 1: Fisher Shut-off on a regression model on six mutually exclusive batches of data. Blue dots represent the overall dataset, while green square dots are the batch or task data. The red line is the model's output after each batch is fed to it. Fisher Shut-off is used from Batch #2 onwards.

182 the shut-off weights must also take part in training. And, as per (1), there is a learning constant
183 required in the regularization to ensure that the right balances between \mathcal{D}_{new} and \mathcal{D}_{old} are met on
184 these weights. This paper provides a novel learning constant determination for this regularization.
185 This is detailed in Appendix A.

186 Also in regression problems, if datasets have batch distributions that are quite far apart from each
187 previous batch, then Fisher-Shutoff may not fully work too. Appendix B shows some of these
188 examples

189 3.2 Fractional Data Retention

190 This is a very simply proposal. The idea is to retain only a fraction of the data trained on the neural
191 network on the earlier tasks or training sets. There is nothing more to this. However, banking on
192 the ideas of selection highlighted in Castro, et al. (2018), whereby data is selected based on their
193 proximity to cluster centers, to be more representative of the classes, this paper uses this as the
194 baseline idea behind its proposal on Fractional Data Retention. Thus in Fractional Data Retention, a
195 fraction of the data within the data cluster is retained and appended in every stage of incremental
196 learning.

197 3.3 Border Control

198 The most important requirement when incrementally learning classes is to ensure that the decision
199 boundaries of the earlier training tasks are protected as much as possible when training on new sets.
200 If data points are used for this purpose, it would seem that, those that lie closest to the decision
201 boundaries after training would be the most important ones to retain, for any succeeding incremental
202 learning tasks. Thus, the Border Control method proposed in this paper exploits this. Ruping (2001)
203 used SVMs to identify these data points as the support vectors that helped define the overall decision
204 boundaries. But with deep learning models or vanilla neural networks, SVM is quite complex and
205 therefore in order to be able retain data points closest to the decision boundaries, a different selection
206 mechanism is required. This selection mechanism could instead be based on selecting data points
207 on how large the absolute errors in sigmoidal outputs are for the applied dataset, as the data points
208 closest to the decision boundaries have this inherent property.

209 Furthermore, since real world data can be quite complex, it would be necessary to not only select data
210 points based on how large their errors in sigmoidal outputs are, but also those points that are far away
211 from them. This is because, given the context of incremental learning where there is a high chance
212 that newer training sets may have data points that could potentially set newer decision boundaries in
213 those farthest regions, these data points would help protect those.

214 Hence in Border Control, the top-k data points that have the largest absolute errors in the sigmoidal
215 outputs and their respective top-k farthest data points are retained in every task or training set
216 for further incremental learning. These points are appended to the newer training sets before further
217 training is applied.

218 4 Sample Results

219 The proposed methods are tested on a toy dataset, as mentioned in Section 1, only to provide some
220 visuals on how the incremental learning progresses using the proposed methods. Figure 2 shows this
221 dataset. A vanilla deep neural network of size, 1000-1000-1000-1000-3, is used for incrementally
222 learning batches of data from this toy dataset. The activation functions for all layers are *ReLU* except
223 for the output which is a *softmax*. All weights are uniformly but randomly initialized with a single
224 random seed to make the results comparable. The weights are also scaled by a $\frac{2}{\sqrt{n}}$ factor to ensure
225 that minimal overfitting occurs during training. Here n is the layer fan-in.

226 To visualize incremental learning on the proposed methods, the dataset is divided into 6 batches
227 with mutually exclusive data points. This gives roughly between 100 to 200 data points in each
228 batch, a size that is commonly used when training neural networks of this size. Figure 3 shows this.
229 In this figure, it can be clearly seen that the batch distributions on the class data for incremental
230 learning do not always form a piece-wise nonlinear fit, and therefore, plain shut-off of weights cannot
231 fully retain knowledge learned earlier. Also, among these distributions, some allowed incremental

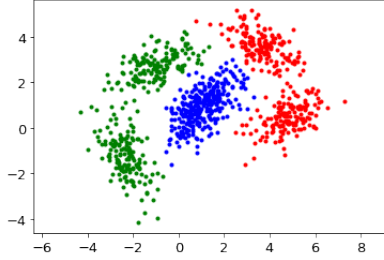


Figure 2: The toy dataset having three nonlinearly arranged classes.

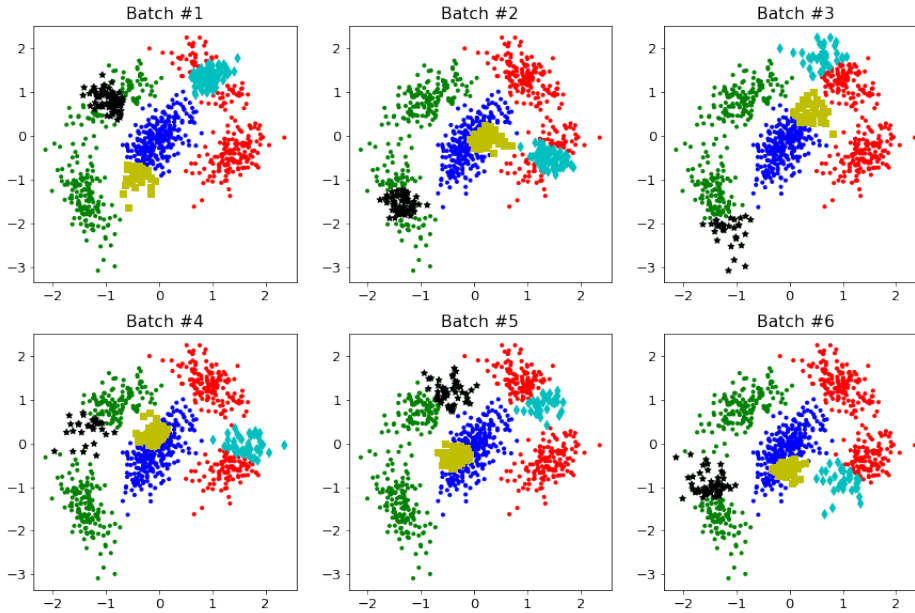


Figure 3: Batches on the toy dataset.

232 learning to happen easily, while others did not, and the distribution shown in Figure 3 is one such.
 233 Table 1 summarizes the results of the proposed methods on this particular distribution. For other batch
 234 distributions, similar results could be achieved. Note here that quite some ML-Ops were required to
 235 achieve the results in Table 1. This was especially the case for those that employed Fisher Shut-off,
 236 since this method has a regularization constant that requires adapting on each batch. Furthermore,
 237 training on each batch was stopped once 100.0% accuracy was obtained on the batch. This left quite
 238 some data points to lie very close to the boundary lines or in some cases just right on them. Thus,
 239 the neural network was very vulnerable to catastrophic forgetting when succeeding batch trainings
 240 occurred as part of incremental learning.

241 However, taking a look at Table 1, it can be seen that when Fisher Shut-off is applied, additional
 242 leverage against catastrophic forgetting occurs on each incremental batch, than when it is not used.
 243 And, among the three methods proposed in this paper, the Border Control method shows much
 244 stronger performance. In Figure 4, sample decision boundaries learned when each incremental batch
 245 is applied to the neural network using Fisher Shut-off and Border Control together is shown. A topk
 246 value of 5 is used for the Border Control. Also, note in Figure 4 that the red circles mark the border
 247 points accumulated on each batch. It can be seen that they clearly assume the data points closest
 248 to the decision boundaries, as well as those far away from it. All with respect to their batches. For
 249 the far away data points, their purpose can be clearly seen between batches #1 and #2, where the
 250 farthest points of class 2 in Batch #1 helped protect the decision boundaries from the data points

Table 1: Performance of proposed methods on the dataset of Figure 2

Method	Sample accuracy on accumulated dataset after Batch ¹					
	#1	#2	#3	#4	#5	#6
No Incremental Learning	100.0%	90.98% ²	92.76%	94.89%	98.32%	96.11%
Fisher Shut-off (FS)	100.0%	99.74%	98.19%	98.88%	98.96%	96.89%
Frac. Data Ret. (FDR)[10%]	100.0%	97.94%	98.39%	98.89%	98.71%	98.67%
FDR[20%]	100.0%	98.71%	98.59%	99.36%	98.97%	98.78%
Border Ctrl. (BC)[topk = 5]	100.0%	100.0%	100.0%	99.84%	100.0%	100.0%
BC[topk = 10]	100.0%	100.0%	100.0%	99.84%	100.0%	100.0%
FS + FDR[10%]	100.0%	99.74%	98.79%	99.52%	99.23%	98.78%
FS + FDR[20%]	100.0%	100.0%	99.19%	99.52%	99.48%	99.11%
FS + BC[topk = 5]	100.0%	100.0%	100.0%	99.84%	100.0%	100.0%
FS + BC[topk = 10]	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%

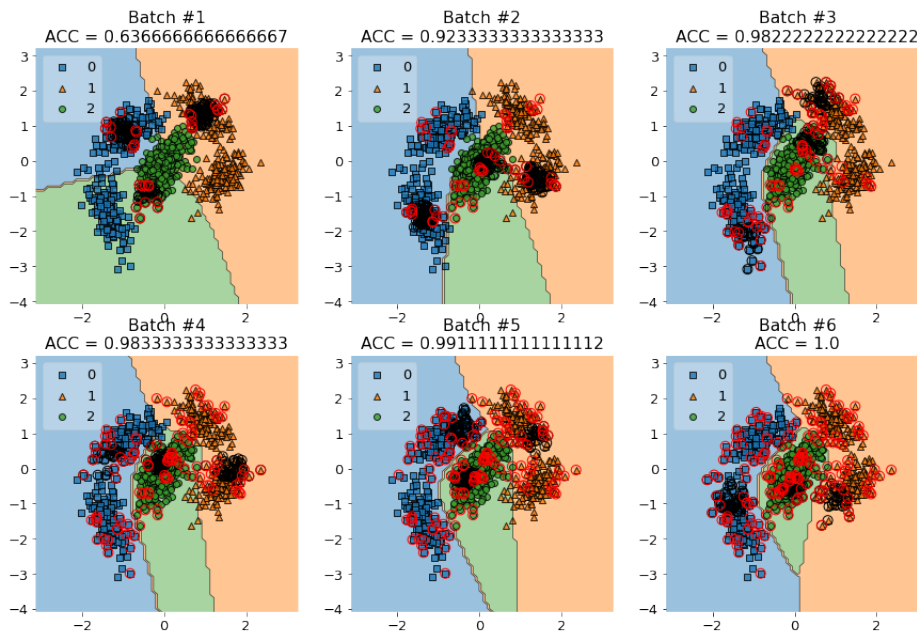


Figure 4: Incremental learning using Fisher Shut-off and Border Control together [$topk = 5$]. Black circles mark the batch data, while the red circles mark the accumulated border points.

251 of class 0 in Batch #2. This means that more complex datasets can be accommodated by simply
 252 applying Border Control. More examples are shown in Appendix C.

253 Also, to add on further to this, for the most difficult incremental learning applications such as learning
 254 new classes as highlighted in Castro, et al. (2018) and He, et al. (2020), Border Control can help
 255 leverage the many issues associated with it like class imbalance, concept drift and so on.

256 5 Conclusion

257 To summarize the work in this paper, three simplified methods for implementing incremental learning
 258 for commercial fine-tuning of pre-trained models was proposed. Results showed that while Border
 259 Control performed the best, Fisher Shut-off was able to leverage the performances. However, dataset
 260 used in this paper was a toy dataset and not one of the benchmark datasets typically used for

¹Incremental learning method applied from Batch #2 onwards.

²Indicates catastrophic forgetting

261 incremental learning. Hence, testing these methods on the benchmark datasets is a potential next step
262 forward. Then, preparing the prerequisites for each model available, like say in `huggingface.com`,
263 for incremental learning must be done so that automation companies or any other AI organization
264 can make use of them. From the methods proposed in this paper, the prerequisites would be: 1)
265 the Shut-off matrix for the neural network weights, and 2) the border points for each of the learned
266 classes. Additionally, an ML-Ops infrastructure can be provided to optimize the performances of
267 the models that employ the Fisher Shut-off method. Metrics like the Backward Transfer (BWT) and
268 Forward Transfer (FWT) proposed in Lopez-Paz & Ranzato (2017) can be used for this purpose.

269 References

- 270 [1] Aich, A. (2021) *Elastic weight consolidation (EWC): Nuts and bolts*. arXiv preprint arXiv:2105.04093.
- 271 [2] Alshalali, T. & Joysula, D. (2018) *Fine-Tuning of Pre-Trained Deep Learning Models with Extreme Learning*
272 *Machine*. IEEE International Conference on Computational Science and Computational Intelligence (CSCI), pp.
273 469-473.
- 274 [3] Castro, F., Marín-Jiménez, M.J., Guil, N., Schmid, C. & Alahari, K. (2018) *End-to-end incremental learning*.
275 Proceedings of the European Conference on Computer Vision (ECCV), pp. 233-248.
- 276 [4] Dif, N. & Elberrichi, Z. (2020) *A New Intra Fine-Tuning Method*. International Journal of Service Science,
277 Management, Engineering, and Technology, **11**(2), pp. 16-40.
- 278 [5] He, J., Mao, R., Shao, Z. & Zhu, F. (2020) *Incremental learning in online scenario*. IEEE/CVF Conference
279 on Computer Vision and Pattern Recognition, pp. 13926-13935.
- 280 [6] Hinton, G., Vinyals, O. & Dean, J. (2015) *Distilling the Knowledge in a Neural Network*. arXiv preprint
281 arXiv:1503.02531, **2**(7).
- 282 [7] Jung, H., Lee, S., Yim, J., Park, S. & Kim, J. (2015) *Joint fine-tuning in deep neural networks for facial*
283 *expression recognition*. Proceedings of the IEEE International Conference on Computer Vision, pp. 2983-2991.
- 284 [8] Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A.A., Milan, K., Quan, J.,
285 Ramalho, T., Grabska-Barwinska, A. & Hassabis, D. (2017) *Overcoming catastrophic forgetting in neural*
286 *networks*. Proceedings of the National Academy of Sciences, **114**(13), pp. 3521-3526.
- 287 [9] Lopez-Paz, D. & Ranzato, M.A. (2017) *Gradient episodic memory for continual learning*. Proceedings of
288 Neural Information Processing Systems (NIPS), pp. 6467-6476.
- 289 [10] Luo, Y., Yin, L., Bai, W. & Mao, K. (2020) *An Appraisal of Incremental Learning Methods*. Entropy,
290 **22**(11), pp. 1190-1216.
- 291 [11] Odena, A., Olah, C. & Shlens, J. (2017) *Conditional image synthesis with auxiliary classifier GANs*.
292 International Conference on Machine Learning (ICML), pp. 2642-2651.
- 293 [12] Polikar, R., Udpa, L., Udpa, S. & Honavar, V. (2001) *Learn++: An Incremental Learning Algorithm for*
294 *Supervised Neural Networks*. IEEE Transactions on Systems, Man, and Cybernetics, part C (applications and
295 reviews), **31**(4), pp. 497-508.
- 296 [13] Qian, X., Zhang, C., Yella, J., Huang, Y., Huang, M.C. & Bom, S. (2021) *Soft sensing model visualization:*
297 *Fine-tuning neural network from what model learned*. IEEE International Conference on Big Data (Big Data),
298 pp. 1900-1908.
- 299 [14] Ren, M., Zeng, W., Yang, B. & Urtasun, R. (2018) *Learning to reweight examples for robust deep learning*.
300 International Conference on Machine Learning (ICML), pp. 4334-4343.
- 301 [15] Ruping, S. (2001) *Incremental learning with support vector machines*. IEEE International Conference on
302 Data Mining, pp. 641-642.
- 303 [16] Rusu, A.A., Rabinowitz, N.C., Desjardins, G., Soyer, H., Kirkpatrick, J., Kavukcuoglu, K., Pascanu, R. &
304 Hadsell, R. (2016) *Progressive neural networks*. arXiv preprint arXiv:1606.04671.
- 305 [17] Too, E., Yujian, L., Njuki, S. & Yingchun, L. (2019) *A comparative study of fine-tuning deep learning*
306 *models for plant disease*. Computers and Electronics in Agriculture, **161**, pp. 272-279.
- 307 [18] Vrbančič, G. & Podgorelec, V. (2020) *Transfer learning with adaptive fine-tuning*. IEEE Access, Volume 8,
308 pp. 196197-196211.
- 309 [19] Wu, Y., Chen, Y.P., Wang, L.J., Ye, Y.C., Liu, Z.C., Guo, Y.D., Zhang, Z.Y. & Fu, Y. (2018) *Incremental*
310 *Classifier Learning with Generative Adversarial Networks*. arXiv preprint arXiv:1802.00853.

311 Checklist

312 The checklist follows the references. Please read the checklist guidelines carefully for information on
313 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
314 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
315 the appropriate section of your paper or providing a brief inline description. For example:

- 316 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 317 • Did you include the license to the code and datasets? **[No]** The code and the data are
318 proprietary.
- 319 • Did you include the license to the code and datasets? **[N/A]**

320 Please do not modify the questions and only use the provided macros for your answers. Note that the
321 Checklist section does not count towards the page limit. In your paper, please delete this instructions
322 block and only keep the Checklist section heading above along with the questions/answers below.

323 1. For all authors...

- 324 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
325 contributions and scope? **[Yes]**
- 326 (b) Did you describe the limitations of your work? **[Yes]** See Sections 4 and 5
- 327 (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
- 328 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
329 them? **[Yes]**

330 2. If you are including theoretical results...

- 331 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]** See Section 3
332 and Appendix A
- 333 (b) Did you include complete proofs of all theoretical results? **[Yes]** See Section 3 and
334 Appendix A

335 3. If you ran experiments...

- 336 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
337 perimental results (either in the supplemental material or as a URL)? **[Yes]** As a
338 supplemental material
- 339 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
340 were chosen)? **[No]** But it’s in the supplemental material
- 341 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
342 ments multiple times)? **[No]** But additional examples are shown in Appendix C
- 343 (d) Did you include the total amount of compute and the type of resources used (e.g., type
344 of GPUs, internal cluster, or cloud provider)? **[N/A]**

345 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 346 (a) If your work uses existing assets, did you cite the creators? **[N/A]**
- 347 (b) Did you mention the license of the assets? **[N/A]**
- 348 (c) Did you include any new assets either in the supplemental material or as a URL? **[N/A]**
349
- 350 (d) Did you discuss whether and how consent was obtained from people whose data you’re
351 using/curating? **[N/A]**
- 352 (e) Did you discuss whether the data you are using/curating contains personally identifiable
353 information or offensive content? **[N/A]**

354 5. If you used crowdsourcing or conducted research with human subjects...

- 355 (a) Did you include the full text of instructions given to participants and screenshots, if
356 applicable? **[N/A]**
- 357 (b) Did you describe any potential participant risks, with links to Institutional Review
358 Board (IRB) approvals, if applicable? **[N/A]**
- 359 (c) Did you include the estimated hourly wage paid to participants and the total amount
360 spent on participant compensation? **[N/A]**

361 **A Appendix**

362 Here, the derivation of regularization constant for Fisher Shut-off method is detailed. To start with,
363 let the neural network be defined as,

$$y = w^T \phi(x, \omega) \quad (7)$$

364 Here, w is the weights of the output layer, $\phi(\cdot)$ is the output of the preceding layer, x is the input and
365 ω represents the rest of the weights of the neural network. Throughout the derivation, we will be
366 dealing with only the output layer, and so the $\phi(x, \omega)$ will be written in short form as Φ from here on.
367 Let e_k denoted the error of fitting in the k -th iteration, and g_k denote the error gradient. This would
368 mean that $g_k = \Phi_k e_k$.

369 Then, given the regularization term in (6), let \tilde{w} denote difference in weights, between the new
370 training and the previous training. This would give the weight updation policy as,

$$w_{k+1} = w_k - \eta g_k - \beta \tilde{w}_k \quad (8)$$

371 If $e_k = w_k^T \Phi_k - Y$, then the error e_{k+1} after the weight updation would equal,

$$\begin{aligned} e_{k+1} &= w_{k+1}^T \Phi_{k+1} - Y \\ &= (w_k - \eta g_k - \beta \tilde{w}_k)^T \Phi_{k+1} - Y \\ &= w_k^T \Phi_{k+1} - \eta g_k^T \Phi_{k+1} - \beta \tilde{w}_k^T \Phi_{k+1} - Y \end{aligned} \quad (9)$$

372 Asuming for simplicity sake that $\Phi_{k+1} \approx \Phi_k + \delta$, then (9) can continue as,

$$\begin{aligned} e_{k+1} &\approx w_k^T \Phi_k - \eta g_k^T \Phi_k - \beta \tilde{w}_k^T \Phi_k - Y + \Delta \\ &\approx e_k - \eta g_k^T \Phi_k - \beta \tilde{w}_k^T \Phi_k \end{aligned} \quad (10)$$

373 Taking the square norm of e_{k+1} in (10), would equate this to,

$$\begin{aligned} \|e_{k+1}\|^2 &= \|e_k\|^2 + \eta^2 \|g_k^T \Phi_k\|^2 + \beta^2 \|\tilde{w}_k^T \Phi_k\|^2 \\ &\quad - 2\eta e_k^T \Phi_k^T g_k - 2\beta e_k^T \Phi_k^T \tilde{w}_k + 2\eta\beta \tilde{w}_k^T \Phi_k \Phi_k^T g_k \end{aligned} \quad (11)$$

374 We require that $\|e_{k+1}\|^2 < \|e_k\|^2$ at all times, so that regularization does not affect the fit at any
375 point during the training. Applying this condition in (11) would give,

$$\eta^2 \|g_k^T \Phi_k\|^2 + \beta^2 \|\tilde{w}_k^T \Phi_k\|^2 - 2\eta e_k^T \Phi_k^T g_k - 2\beta e_k^T \Phi_k^T \tilde{w}_k + 2\eta\beta \tilde{w}_k^T \Phi_k \Phi_k^T g_k < 0 \quad (12)$$

376 Then, taking the partial derivatives of (12) w.r.t η and β would give the following equations to be
377 satisfied:

$$\eta \|g_k^T \Phi_k\|^2 - e_k^T \Phi_k^T g_k + \beta \tilde{w}_k^T \Phi_k \Phi_k^T g_k = 0 \quad (13)$$

$$\beta \|\tilde{w}_k^T \Phi_k\|^2 - e_k^T \Phi_k^T \tilde{w}_k + \eta \tilde{w}_k^T \Phi_k \Phi_k^T g_k = 0 \quad (14)$$

378 Solving, (13) and (14) can result in negative η and β , which is not acceptable, and so to simplify the
 379 solution, we neglect the β -term in (13). Doing so we get,

$$\eta = \frac{e_k^T \Phi_k^T g_k}{\|g_k^T \Phi_k\|^2}, \quad (15)$$

$$\beta = \frac{e_k^T \Phi_k^T \tilde{w}_k - \eta \tilde{w}_k^T \Phi_k \Phi_k^T g_k}{\|\tilde{w}_k^T \Phi_k\|^2}$$

380 Then, substituting for g_k and opening up the norms, we get,

$$\eta = \frac{e_k^T \Phi_k^T \Phi_k e_k}{e_k^T (\Phi_k^T \Phi_k) (\Phi_k^T \Phi_k) e_k}, \quad (16)$$

$$\beta = \frac{e_k^T \Phi_k^T \tilde{w}_k - \eta \tilde{w}_k^T (\Phi_k \Phi_k^T) \Phi_k e_k}{\tilde{w}_k^T (\Phi_k \Phi_k^T) \tilde{w}_k}$$

381 Simplifying (16) gives,

$$\eta = \frac{e_k^T e_k}{e_k^T \Phi_k^T \Phi_k e_k}, \quad (17)$$

$$\beta = (1 - \eta) \frac{\tilde{w}_k^T \Phi_k e_k}{\tilde{w}_k^T \tilde{w}_k}$$

382 Equation (17) gives the raw form for both η and β to be regulated. However, this will be further
 383 simplified for computing purposes, but will be used as a basis.

384 Since the errors e_k get smaller as the neural network fits the data, using them in learning constants
 385 will only slow down the fits. A common way to overcome this is by replacing e_k with all ones.
 386 Similarly, for the \tilde{w}_k , all weights that are to be regularized are replaced with ones. If we denote
 387 the weights to be regularized as w_r , and there are m patterns in the dataset with n weights to be
 388 regularized, the η and β computations become,

$$\eta = \frac{1}{\|\Phi_k\|^2}, \quad (18)$$

$$\beta = \frac{\alpha}{mn} \sum_{i:w \in w_r} \Phi_{i,k}$$

389 Here, α represents the $(1 - \eta)$ -term in (17). This constant will not necessarily take the computed η
 390 when being regulated. Instead, this constant will have to be adapted each time for every incremental
 391 batch applied to the neural network.

392 The reason why the computed η is not used for the α adaptation is because this η can sometimes
 393 become too small in the adaptation, that the $1 - \eta$ would always tend towards 1. When this was
 394 empirically tested on the toy dataset, the regularization was found at times to have gone too strong
 395 that the fit never happened. ML-Ops on the α found that this constant is not always 1, and can be
 396 anywhere between 0 and 1, or higher in some cases.

397 **B Appendix**

398 Additional examples on the regression problem with Fisher Shut-off. Fisher Shut-off could not be used completely, and regularization had to take over for some batches. Figures 5 and 6 show this.

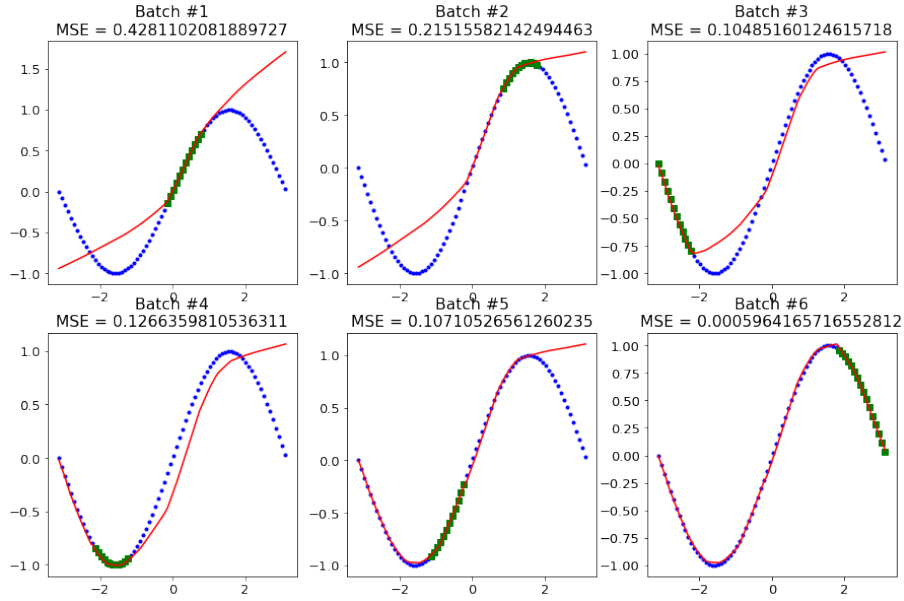


Figure 5: Complete Fisher Shut-off is used in Batches #2 and #6. Batches #3, #4 and #5 are regularized.

399

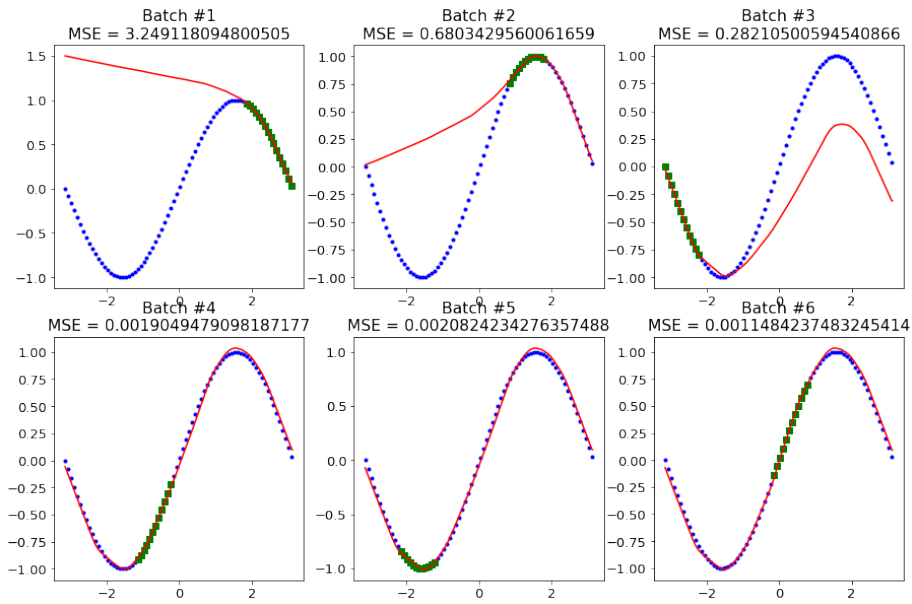


Figure 6: Complete Fisher Shut-off is used in Batches #2 and #6. Batches #3 and #4 are regularized. Batch #5 is fine-tuned

400 **C Appendix**

401 In this appendix, additional examples on the toy dataset classification is shown. Figures 7 and 8 show
 402 the batch distributions considered. Among these, Figure 8 has more cases in which the farthest points
 403 in Border Control can play a vital role in retaining previously learned knowledge. Tables 2 and 3
 summarize their performances.

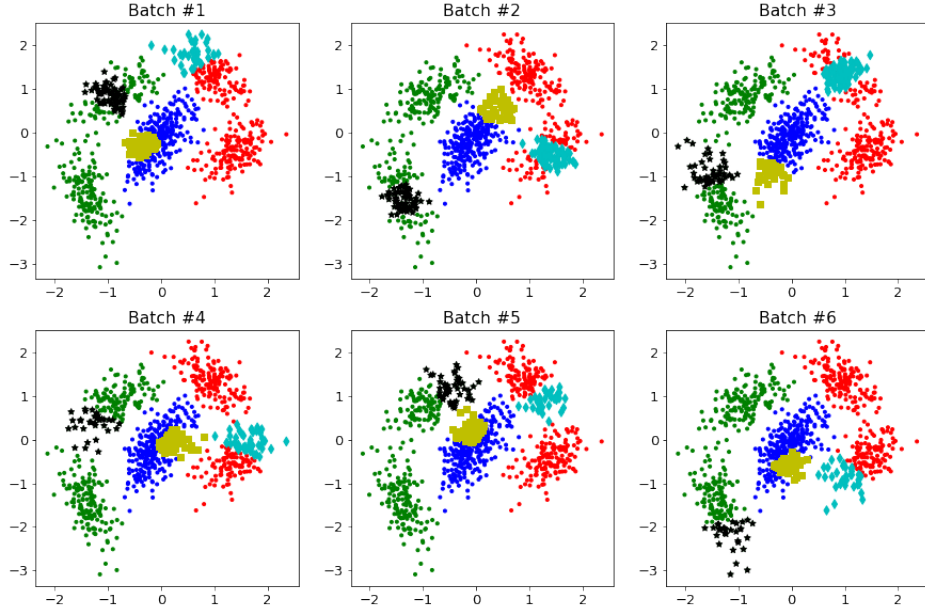


Figure 7: Another batch distribution on the toy dataset.

404

Table 2: Performance of proposed methods on the dataset of Figure 7

Method	Sample accuracy on accumulated dataset after Batch					
	#1	#2	#3	#4	#5	#6
No Incremental Learning	100.0%	73.14%	92.07%	97.56%	88.33%	91.67%
Fisher Shut-off (FS)	100.0%	99.43%	98.26%	99.54%	92.85%	97.78%
Frac. Data Ret. (FDR)[10%]	100.0%	97.43%	96.13%	98.93%	98.75%	96.78%
FDR[20%]	100.0%	98.86%	98.84%	99.69%	99.75%	98.89%
Border Ctrl. (BC)[topk = 5]	100.0%	100.0%	100.0%	100.0%	100.0%	99.78%
BC[topk = 10]	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%
FS + FDR[10%]	100.0%	99.71%	98.84%	99.54%	98.75%	98.33%
FS + FDR[20%]	100.0%	100.0%	99.23%	99.85%	99.87%	98.89%
FS + BC[topk = 5]	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%
FS + BC[topk = 10]	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%

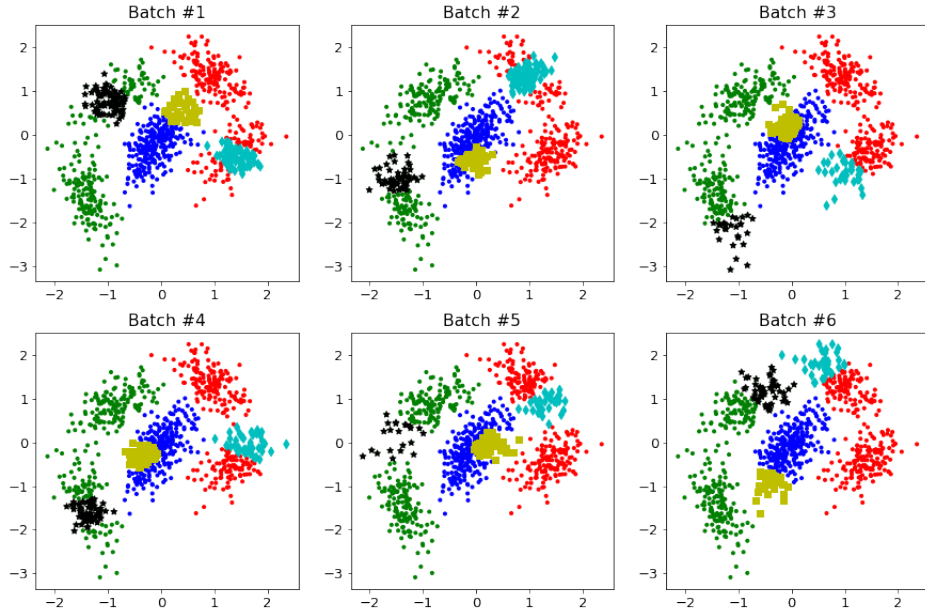


Figure 8: Yet another batch distribution on the toy dataset.

Table 3: Performance of proposed methods on the dataset of Figure 8

Method	Sample accuracy on accumulated dataset after Batch					
	#1	#2	#3	#4	#5	#6
No Incremental Learning	100.0%	89.34%	93.80%	95.60%	97.81%	84.78%
Fisher Shut-off (FS)	100.0%	95.36%	99.59%	99.55%	99.36%	94.56%
Frac. Data Ret. (FDR)[10%]	100.0%	95.08%	96.07%	97.42%	99.61%	98.67%
FDR[20%]	100.0%	98.36%	98.97%	99.85%	100.0%	99.67%
Border Ctrl. (BC)[topk = 5]	100.0%	100.0%	99.79%	99.85%	100.0%	100.0%
BC[topk = 10]	100.0%	100.0%	100.0%	100.0%	99.74%	100.0%
FS + FDR[10%]	100.0%	95.36%	99.79%	98.48%	99.61%	98.89%
FS + FDR[20%]	100.0%	98.36%	99.79%	99.69%	100.0%	99.67%
FS + BC[topk = 5]	100.0%	100.0%	100.0%	100.0%	99.74%	100.0%
FS + BC[topk = 10]	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%