

---

# Sequential Attention for Feature Selection

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Feature selection is the problem of selecting a subset of features for a machine learn-  
2 ing model that maximizes model quality subject to a budget constraint. For neural  
3 networks, prior methods, including those based on  $\ell_1$  regularization, attention,  
4 and other techniques, typically select the entire feature subset in one evaluation  
5 round, ignoring the residual value of features during selection, i.e., the marginal  
6 contribution of a feature given that other features have already been selected. We  
7 propose a feature selection algorithm called Sequential Attention that achieves  
8 state-of-the-art empirical results for neural networks. This algorithm is based on an  
9 efficient one-pass implementation of greedy forward selection and uses attention  
10 weights at each step as a proxy for feature importance. We give theoretical insights  
11 into our algorithm for linear regression by showing that an adaptation to this setting  
12 is equivalent to the classical Orthogonal Matching Pursuit (OMP) algorithm, and  
13 thus inherits all of its provable guarantees. Our theoretical and empirical analyses  
14 offer new explanations towards the effectiveness of attention and its connections to  
15 overparameterization, which may be of independent interest.

## 16 1 Introduction

17 Feature selection is a classic problem in machine learning and statistics where one is asked to find  
18 a subset of  $k$  features from a larger set of  $d$  features, such that the prediction quality of the model  
19 trained using the subset of features is maximized. Finding a small and high-quality feature subset is  
20 desirable for many reasons: improving model interpretability, reducing inference latency, decreasing  
21 model size, regularization, and removing redundant or noisy features to improve generalization. We  
22 direct the reader to [Li et al. \(2017b\)](#) for a survey on the role of feature selection in machine learning.

23 The widespread success of deep learning has prompted an intense study of feature selection algorithms  
24 for neural networks, especially in the supervised setting. While many methods have been proposed,  
25 we focus on a line of work that studies the use of *attention for feature selection*. The attention  
26 mechanism in machine learning roughly refers to applying a trainable softmax mask to a given layer.  
27 This allows the model to “focus” on certain important signals during training. Attention has recently  
28 led to major breakthroughs in computer vision, natural language processing, and several other areas  
29 of machine learning ([Vaswani et al., 2017](#)). For feature selection, the works of [Wang et al. \(2014\)](#);  
30 [Gui et al. \(2019\)](#); [Skrlj et al. \(2020\)](#); [Wojtas and Chen \(2020\)](#); [Liao et al. \(2021\)](#) all present new  
31 approaches for feature attribution, ranking, and selection that are inspired by attention.

32 One problem with naively using attention for feature selection is that it can ignore the *residual values*  
33 of features, i.e., the marginal contribution a feature has on the loss conditioned on previously-selected  
34 features being in the model. This can lead to several problems such as selecting redundant features or  
35 ignoring features that are uninformative in isolation but valuable in the presence of others.

36 This work introduces the *Sequential Attention* algorithm for supervised feature selection. Our algo-  
37 rithm addresses the shortcomings above by using attention-based selection *adaptively* over multiple

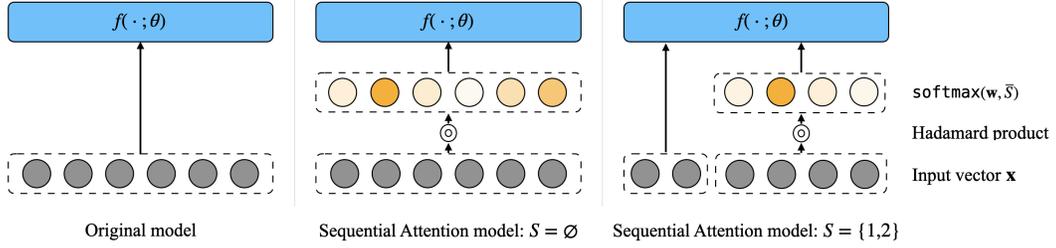


Figure 1: Sequential attention applied to model  $f(\cdot; \theta)$ . At each step, the selected features  $i \in S$  are used as direct inputs to the model and the unselected features  $i \notin S$  are downscaled by the scalar value  $\text{softmax}_i(\mathbf{w}, \bar{S})$ , where  $\mathbf{w} \in \mathbb{R}^d$  is the vector of learned attention weights and  $\bar{S} = [d] \setminus S$ .

38 rounds. Further, Sequential Attention simplifies earlier attention-based approaches by directly training  
 39 one global feature mask instead of aggregating many instance-wise feature masks. This technique  
 40 reduces the overhead of our algorithm, eliminates the toil of tuning unnecessary hyperparameters,  
 41 *works directly with any differentiable model architecture*, and offers an efficient streaming implemen-  
 42 tation. Empirically, Sequential Attention achieves state-of-the-art feature selection results for neural  
 43 networks on standard benchmarks. The code for our algorithm and experiments is publicly available.<sup>1</sup>

44 **Sequential Attention.** Our starting point for Sequential Attention is the well-known greedy forward  
 45 selection algorithm, which repeatedly selects the feature with the *largest marginal improvement* in  
 46 model loss when added to the set of currently selected features (see, e.g., [Das and Kempe \(2011\)](#)  
 47 and [Elenberg et al. \(2018\)](#)). Greedy forward selection is known to select high-quality features, but  
 48 requires training  $O(kd)$  models and is thus impractical for many modern machine learning problems.  
 49 To reduce this cost, one natural idea is to only train  $k$  models, where the model trained in each step  
 50 approximates the marginal gains of all  $O(d)$  unselected features. Said another way, we can relax  
 51 the greedy algorithm to fractionally consider all  $O(d)$  feature candidates simultaneously rather than  
 52 computing their exact marginal gains one-by-one with separate models. We implement this idea by  
 53 introducing a new set of trainable variables  $\mathbf{w} \in \mathbb{R}^d$  that represent *feature importance*, or *attention*  
 54 *logits*. In each step, we select the feature with maximum importance and add it to the selected set. To  
 55 ensure the score-augmented models (1) have differentiable architectures and (2) are encouraged to  
 56 hone in on the best unselected feature, we take the *softmax* of the importance scores and multiply  
 57 each input feature value by its corresponding softmax value as illustrated in Figure 1.

58 Formally, given a dataset  $\mathbf{X} \in \mathbb{R}^{n \times d}$  represented as a matrix with  $n$  rows of examples and  $d$  feature  
 59 columns, suppose we want to select  $k$  features. Let  $f(\cdot; \theta)$  be a differentiable model, e.g., a neural  
 60 network, that outputs the predictions  $f(\mathbf{X}; \theta)$ . Let  $\mathbf{y} \in \mathbb{R}^n$  be the labels,  $\ell(f(\mathbf{X}; \theta), \mathbf{y})$  be the loss  
 61 between the model’s predictions and the labels, and  $\circ$  be the Hadamard product. Sequential Attention  
 62 outputs a subset  $S \subseteq [d] := \{1, 2, \dots, d\}$  of  $k$  feature indices, and is presented below in Algorithm 1.

---

**Algorithm 1** Sequential Attention for feature selection.

---

- 1: **function** SEQUENTIALATTENTION(dataset  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , labels  $\mathbf{y} \in \mathbb{R}^n$ , model  $f$ , loss  $\ell$ , size  $k$ )
- 2:   Initialize  $S \leftarrow \emptyset$
- 3:   **for**  $t = 1$  to  $k$  **do**
- 4:     Let  $(\theta^*, \mathbf{w}^*) \leftarrow \arg \min_{\theta, \mathbf{w}} \ell(f(\mathbf{X} \circ \mathbf{W}; \theta), \mathbf{y})$ , where  $\mathbf{W} = \mathbf{1}_n \text{softmax}(\mathbf{w}, \bar{S})^\top$  for

$$\text{softmax}_i(\mathbf{w}, \bar{S}) := \begin{cases} 1 & \text{if } i \in S \\ \frac{\exp(\mathbf{w}_i)}{\sum_{j \in \bar{S}} \exp(\mathbf{w}_j)} & \text{if } i \in \bar{S} := [d] \setminus S \end{cases} \quad (1)$$

- 5:     Set  $i^* \leftarrow \arg \max_{i \notin S} \mathbf{w}_i^*$    ▷ unselected feature with largest attention weight
  - 6:     Update  $S \leftarrow S \cup \{i^*\}$
  - 7:   **return**  $S$
- 

<sup>1</sup>Anonymous while under review

63 **Theoretical guarantees.** We give provable guarantees for Sequential Attention for least squares  
 64 linear regression by analyzing a variant of the algorithm called *regularized linear Sequential Attention*.  
 65 This variant (1) uses Hadamard product overparameterization directly between the attention weights  
 66 and feature values without normalizing the attention weights via  $\text{softmax}(\mathbf{w}, \bar{S})$ , and (2) adds  $\ell_2$   
 67 regularization to the objective, hence the “linear” and “regularized” terms. Note that  $\ell_2$  regularization,  
 68 or *weight decay*, is common practice when using gradient-based optimizers (Tibshirani, 2021). We  
 69 give theoretical and empirical evidence that replacing the softmax by different overparameterization  
 70 schemes leads to similar results (Section 4.2) while offering more tractable analysis. In particular, our  
 71 main result shows that regularized linear Sequential Attention has the same provable guarantees as  
 72 the celebrated *Orthogonal Matching Pursuit* (OMP) algorithm of Pati et al. (1993) for sparse linear  
 73 regression, without making any assumptions on the design matrix or response vector.

74 **Theorem 1.1.** *For linear regression, regularized linear Sequential Attention is equivalent to OMP.*

75 We prove this equivalence using a novel two-step argument. First, we show that regularized linear  
 76 Sequential Attention is equivalent to a greedy version of LASSO (Tibshirani, 1996), which Luo  
 77 and Chen (2014) call *Sequential LASSO*. Prior to our work, however, Sequential LASSO was only  
 78 analyzed in a restricted “sparse signal plus noise” setting, offering limited insight into its success in  
 79 practice. Second, we prove that Sequential LASSO is equivalent to OMP in the fully general setting  
 80 for linear regression by analyzing the geometry of the associated polyhedra. This ultimately allows  
 81 us to transfer the guarantees of OMP to Sequential Attention.

82 **Theorem 1.2.** *For linear regression, Sequential LASSO (Luo and Chen, 2014) is equivalent to OMP.*

83 We present the full argument for our results in Section 3. This analysis takes significant steps towards  
 84 explaining the success of attention in feature selection and the various theoretical phenomena at play.

85 **Towards understanding attention.** An important property of OMP is that it provably approximates  
 86 the marginal gains of features—Das and Kempe (2011) showed that for any subset of features, the  
 87 gradient of the least squares loss at its sparse minimizer approximates the marginal gains up to a factor  
 88 that depends on the *sparse condition numbers* of the design matrix. This suggests that Sequential  
 89 Attention could also approximate some notion of the marginal gains for more sophisticated models  
 90 when selecting the next-best feature. We observe this phenomenon empirically in our marginal gain  
 91 experiments in Appendix B.6. These results also help refine the widely-assumed conjecture that  
 92 attention weights correlate with feature importances by specifying an exact measure of “importance”  
 93 at play. Since a countless number of feature importance definitions are used in practice, it is important  
 94 to understand which best explains how the attention mechanism works.

95 **Connections to overparameterization.** In our analysis of regularized linear Sequential Attention  
 96 for linear regression, we do not use the presence of the softmax in the attention mechanism—rather,  
 97 the crucial ingredient in our analysis is the Hadamard product parameterization of the learned weights.  
 98 We conjecture that the empirical success of attention-based feature selection is primarily due to the  
 99 explicit overparameterization.<sup>2</sup> Indeed, our experiments in Section 4.2 verify this claim by showing  
 100 that if we substitute the softmax in Sequential Attention with a number of different (normalized)  
 101 overparameterized expressions, we achieve nearly identical performance. This line of reasoning is also  
 102 supported in the recent work of Ye et al. (2021), who claim that attention largely owes its success to  
 103 the “smoother and stable [loss] landscapes” induced by Hadamard product overparameterization.

## 104 1.1 Related work

105 Here we discuss recent advances in supervised feature selection for deep neural networks (DNNs)  
 106 that are the most related to our empirical results. In particular, we omit a discussion of a large body of  
 107 works on unsupervised feature selection (Zou et al., 2015; Altschuler et al., 2016; Balm et al., 2019).

108 The *group LASSO* method has been applied to DNNs to achieve structured sparsity by pruning  
 109 neurons (Alvarez and Salzmann, 2016) and even filters or channels in convolutional neural net-  
 110 works (Lebedev and Lempitsky, 2016; Wen et al., 2016; Li et al., 2017a). It has also be applied for  
 111 feature selection (Zhao et al., 2015; Li et al., 2016; Scardapane et al., 2017; Lemhadri et al., 2021).

<sup>2</sup>Note that overparameterization here refers to the addition of  $d$  trainable variables in the Hadamard product overparameterization, not the other use of the term that refers to the use of a massive number of parameters in neural networks, e.g., in Bubeck and Sellke (2021).

112 While the LASSO is the most widely-used method for relaxing the  $\ell_0$  sparsity constraint in feature  
 113 selection, several recent works have proposed new relaxations based on *stochastic gates* (Srinivas  
 114 et al., 2017; Louizos et al., 2018; Baln et al., 2019; Trelin and Procházka, 2020; Yamada et al., 2020).  
 115 This approach introduces (learnable) Bernoulli random variables for each feature during training, and  
 116 minimizes the expected loss over realizations of the 0-1 variables (accepting or rejecting features).

117 There are several other recent approaches for DNN feature selection. Roy et al. (2015) explore using  
 118 the magnitudes of weights in the first hidden layer to select features. Lu et al. (2018) designed the  
 119 DeepPINK architecture, extending the idea of *knockoffs* (Benjamini et al., 2001) to neural networks.  
 120 Here, each feature competes with a “knockoff” version of the original feature; if the knockoff wins,  
 121 the feature is removed. Borisov et al. (2019) introduced the *CancelOut* layer, which suppresses  
 122 irrelevant features via independent per-feature activation functions that act as (soft) bitmasks.

123 In contrast to these differentiable approaches, the combinatorial optimization literature is rich with  
 124 greedy algorithms that have applications in machine learning (Zadeh et al., 2017; Fahrback et al.,  
 125 2019b,a; Chen et al., 2021; Halabi et al., 2022; Bilmes, 2022). In fact, most influential feature selection  
 126 algorithms from this literature are sequential, e.g., greedy forward and backward selection (Ye and  
 127 Sun, 2018; Das et al., 2022), Orthogonal Matching Pursuit (Pati et al., 1993), and several information-  
 128 theoretic methods (Fleuret, 2004; Ding and Peng, 2005; Bennasar et al., 2015). These approaches,  
 129 however, are not normally tailored to neural networks, and can suffer from quality, efficiency, or both.

130 Lastly, this paper studies *global* feature selection, i.e., selecting the same subset of features across  
 131 all training examples, whereas many works consider *local* (or instance-wise) feature selection. This  
 132 problem is more related to model interpretability, and is better known as *feature attribution* or  
 133 *saliency maps*. These methods naturally lead to global feature selection methods by aggregating their  
 134 instance-wise scores (Cancela et al., 2020). Instance-wise feature selection has been explored using a  
 135 variety of techniques, including gradients (Smilkov et al., 2017; Sundararajan et al., 2017; Srinivas  
 136 and Fleuret, 2019), attention (Arik and Pfister, 2021; Ye et al., 2021), mutual information (Chen et al.,  
 137 2018), and Shapley values from cooperative game theory (Lundberg and Lee, 2017).

## 138 2 Preliminaries

139 Before discussing our theoretical guarantees for Sequential Attention in Section 3, we present several  
 140 known results about feature selection for linear regression, also called *sparse linear regression*. Recall  
 141 that in the least squares linear regression problem, we have

$$\ell(f(\mathbf{X}; \boldsymbol{\theta}), \mathbf{y}) = \|f(\mathbf{X}; \boldsymbol{\theta}) - \mathbf{y}\|_2^2 = \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2. \quad (2)$$

142 We work in the most challenging setting for obtaining relative error guarantees for this objective by  
 143 making *no distributional assumptions* on  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , i.e., we seek  $\tilde{\boldsymbol{\theta}} \in \mathbb{R}^d$  such that

$$\|\mathbf{X}\tilde{\boldsymbol{\theta}} - \mathbf{y}\|_2^2 \leq \kappa \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2, \quad (3)$$

144 for some  $\kappa = \kappa(\mathbf{X}) > 0$ , where  $\mathbf{X}$  is not assumed to follow any particular input distribution. This  
 145 is far more applicable in practice than assuming the entries of  $\mathbf{X}$  are i.i.d. Gaussian. In large-scale  
 146 applications, the number of examples  $n$  often greatly exceeds the number of features  $d$ , resulting in  
 147 an optimal loss that is nonzero. Thus, we focus on the *overdetermined* regime and refer to Price et al.  
 148 (2022) for an excellent discussion on the long history of this problem.

149 **Notation.** Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  be the design matrix with  $\ell_2$  unit columns and let  $\mathbf{y} \in \mathbb{R}^n$  be the  
 150 response vector, also assumed to be an  $\ell_2$  unit vector.<sup>3</sup> For  $S \subseteq [d]$ , let  $\mathbf{X}_S$  denote the  $n \times |S|$  matrix  
 151 consisting of the columns of  $\mathbf{X}$  indexed by  $S$ . For singleton sets  $S = \{j\}$ , we write  $\mathbf{X}_j$  for  $\mathbf{X}_{\{j\}}$ .  
 152 Let  $\mathbf{P}_S := \mathbf{X}_S \mathbf{X}_S^+$  denote the projection matrix onto the column span  $\text{colspan}(\mathbf{X}_S)$  of  $\mathbf{X}_S$ , where  
 153  $\mathbf{X}_S^+$  denotes the pseudoinverse of  $\mathbf{X}_S$ . Let  $\mathbf{P}_S^\perp = \mathbf{I}_n - \mathbf{P}_S$  denote the projection matrix onto the  
 154 orthogonal complement of  $\text{colspan}(\mathbf{X}_S)$ .

155 **Feature selection algorithms for linear regression.** Perhaps the most natural algorithm for sparse  
 156 linear regression is greedy forward selection, which was shown to have guarantees of the form of (3) in

<sup>3</sup>These assumptions are without loss of generality by scaling.

157 the breakthrough works of [Das and Kempe \(2011\)](#); [Elenberg et al. \(2018\)](#), where  $\kappa = \kappa(\mathbf{X})$  depends  
 158 on *sparse condition numbers* of  $\mathbf{X}$ , i.e., the spectrum of  $\mathbf{X}$  restricted to a subset of its columns. Greedy  
 159 forward selection can be expensive in practice, but these works also prove analogous guarantees for  
 160 the more efficient Orthogonal Matching Pursuit algorithm, which we present formally in Algorithm 2.

---

**Algorithm 2** Orthogonal Matching Pursuit ([Pati et al., 1993](#)).

---

1: **function** OMP(design matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , response  $\mathbf{y} \in \mathbb{R}^n$ , size constraint  $k$ )  
 2:   Initialize  $S \leftarrow \emptyset$   
 3:   **for**  $t = 1$  to  $k$  **do**  
 4:     Set  $\beta_S^* \leftarrow \arg \min_{\beta \in \mathbb{R}^S} \|\mathbf{X}_S \beta - \mathbf{y}\|_2^2$   
 5:     Let  $i^* \notin S$  maximize ▷ maximum correlation with residual  
       
$$\langle \mathbf{X}_{i^*}, \mathbf{y} - \mathbf{X}_S \beta_S^* \rangle^2 = \langle \mathbf{X}_{i^*}, \mathbf{y} - \mathbf{P}_S \mathbf{y} \rangle^2 = \langle \mathbf{X}_{i^*}, \mathbf{P}_S^\perp \mathbf{y} \rangle^2$$
  
 6:     Update  $S \leftarrow S \cup \{i^*\}$   
 7:   **return**  $S$

---

161 The LASSO algorithm ([Tibshirani, 1996](#)) is another popular feature selection method, which simply  
 162 adds  $\ell_1$ -regularization to the objective in Equation (2). Theoretical guarantees for LASSO are known  
 163 in the *underdetermined* regime ([Donoho and Elad, 2003](#); [Candes and Tao, 2006](#)), but it is an open  
 164 problem whether LASSO has the guarantees of Equation (3). Sequential LASSO is a related algorithm  
 165 that uses LASSO to select features one by one. [Luo and Chen \(2014\)](#) analyzed this algorithm in  
 166 a specific parameter regime, but until our work, *no relative error guarantees were known in full*  
 167 *generality (e.g., the overdetermined regime)*. We present the Sequential LASSO in Algorithm 3.

---

**Algorithm 3** Sequential LASSO ([Luo and Chen, 2014](#)).

---

1: **function** SEQUENTIALLASSO(design matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , response  $\mathbf{y} \in \mathbb{R}^n$ , size constraint  $k$ )  
 2:   Initialize  $S \leftarrow \emptyset$   
 3:   **for**  $t = 1$  to  $k$  **do**  
 4:     Let  $\beta^*(\lambda, S)$  denote the optimal solution to  
       
$$\arg \min_{\beta \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{X} \beta - \mathbf{y}\|_2^2 + \lambda \|\beta_{\overline{S}}\|_1 \tag{4}$$
  
 5:     Set  $\lambda^*(S) \leftarrow \sup\{\lambda > 0 : \beta^*(\lambda, S)_{\overline{S}} \neq \mathbf{0}\}$  ▷ largest  $\lambda$  with nonzero LASSO on  $\overline{S}$   
 6:     Let  $A(S) = \lim_{\varepsilon \rightarrow 0} \{i \in \overline{S} : \beta^*(\lambda^* - \varepsilon, S)_i \neq 0\}$   
 7:     Select any  $i^* \in A(S)$  ▷ non-empty by Lemma 3.5  
 8:     Update  $S \leftarrow S \cup \{i^*\}$   
 9:   **return**  $S$

---

168 Note that Sequential LASSO as stated requires a search for the optimal  $\lambda^*$  in each step. In practice,  $\lambda$   
 169 can simply be set to a large enough value to obtain similar results, since beyond a critical value of  $\lambda$ ,  
 170 the feature ranking according to LASSO coefficients does not change ([Efron et al., 2004](#)).

### 171 3 Equivalence for least squares: OMP and Sequential Attention

172 In this section, we show that the following algorithms are equivalent for least squares linear regression:  
 173 regularized linear Sequential Attention, Sequential LASSO, and Orthogonal Matching Pursuit.

#### 174 3.1 Regularized linear Sequential Attention and Sequential LASSO

175 We start by formalizing a modification to Sequential Attention that admits provable guarantees.

176 **Definition 3.1** (Regularized linear Sequential Attention). Let  $S \subseteq [d]$  be the set of currently se-  
 177 lected features. We define the *regularized linear Sequential Attention* objective by removing the  
 178  $\text{softmax}(\mathbf{w}, \overline{S})$  normalization in Algorithm 1 and introducing  $\ell_2$  regularization on the importance

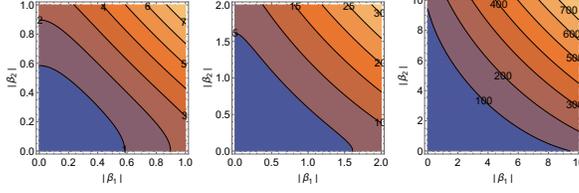


Figure 2: Contour plot of  $Q^*(\beta \circ \beta)$  for  $\beta \in \mathbb{R}^2$  at different zoom-levels of  $|\beta_i|$ .

179 weights  $\mathbf{w} \in \mathbb{R}^{\bar{S}}$  and model parameters  $\theta \in \mathbb{R}^d$  restricted to  $\bar{S}$ . That is, we consider the objective

$$\min_{\mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}^d} \|\mathbf{X}(\mathbf{s}(\mathbf{w}) \circ \theta) - \mathbf{y}\|_2^2 + \frac{\lambda}{2} (\|\mathbf{w}\|_2^2 + \|\theta_{\bar{S}}\|_2^2), \quad (5)$$

180 where  $\mathbf{s}(\mathbf{w}) \circ \theta$  denotes the Hadamard product,  $\theta_{\bar{S}} \in \mathbb{R}^{\bar{S}}$  is  $\theta$  restricted to indices in  $\bar{S}$ , and

$$\mathbf{s}_i(\mathbf{w}, \bar{S}) := \begin{cases} 1 & \text{if } i \in S, \\ \mathbf{w}_i & \text{if } i \notin S. \end{cases}$$

181 By a simple argument due to Hoff (2017), the objective function in (5) is equivalent to

$$\min_{\theta \in \mathbb{R}^d} \|\mathbf{X}\theta - \mathbf{y}\|_2^2 + \lambda \|\theta_{\bar{S}}\|_1. \quad (6)$$

182 It follows that attention (or more generally overparameterization by trainable weights  $\mathbf{w}$ ) can be seen  
 183 as a way to implement  $\ell_1$  regularization for least squares linear regression, i.e., the LASSO (Tibshirani,  
 184 1996). This connection between overparameterization and  $\ell_1$  regularization has also been observed in  
 185 several other recent works (Vaskevicius et al., 2019; Zhao et al., 2022; Tibshirani, 2021).

186 By this transformation and reasoning, regularized linear Sequential Attention can be seen as iteratively  
 187 using the LASSO with  $\ell_1$  regularization applied only to the unselected features—which is precisely  
 188 the Sequential LASSO algorithm in Luo and Chen (2014). If we instead use  $\text{softmax}(\mathbf{w}, \bar{S})$  as in (1),  
 189 then this only changes the choice of regularization, as shown in Lemma 3.2 (proof in Appendix A.3).

190 **Lemma 3.2.** Let  $D : \mathbb{R}^d \rightarrow \mathbb{R}^{\bar{S}}$  be the function defined by  $D(\mathbf{w})_i = 1/\text{softmax}_i^2(\mathbf{w}, \bar{S})$ , for  $i \in \bar{S}$ .  
 191 Denote its range and preimage by  $\text{ran}(D) \subseteq \mathbb{R}^{\bar{S}}$  and  $D^{-1}(\cdot) \subseteq \mathbb{R}^d$ , respectively. Moreover, define  
 192 the functions  $Q : \text{ran}(D) \rightarrow \mathbb{R}$  and  $Q^* : \mathbb{R}^{\bar{S}} \rightarrow \mathbb{R}$  by

$$Q(\mathbf{q}) = \inf_{\mathbf{w} \in D^{-1}(\mathbf{q})} \|\mathbf{w}\|_2^2 \quad \text{and} \quad Q^*(\mathbf{x}) = \inf_{\mathbf{q} \in \text{ran}(D)} \left( \sum_{i \in \bar{S}} \mathbf{x}_i \mathbf{q}_i + Q(\mathbf{q}) \right).$$

193 Then, the following two optimization problems with respect to  $\beta \in \mathbb{R}^d$  are equivalent:

$$\inf_{\substack{\beta \in \mathbb{R}^d \\ \text{s.t. } \beta = \text{softmax}(\mathbf{w}, \bar{S}) \circ \theta \\ \mathbf{w} \in \mathbb{R}^d, \theta \in \mathbb{R}^d}} \|\mathbf{X}\beta - \mathbf{y}\|_2^2 + \frac{\lambda}{2} (\|\mathbf{w}\|_2^2 + \|\theta_{\bar{S}}\|_2^2) = \inf_{\beta \in \mathbb{R}^d} \|\mathbf{X}\beta - \mathbf{y}\|_2^2 + \frac{\lambda}{2} Q^*(\beta \circ \beta). \quad (7)$$

194 We present contour plots of  $Q^*(\beta \circ \beta)$  for  $\beta \in \mathbb{R}^2$  in Figure 2. These plots suggest that  $Q^*(\beta \circ \beta)$   
 195 is a concave regularizer when  $|\beta_1| + |\beta_2| > 2$ , which would thus approximate the  $\ell_0$  regularizer and  
 196 induce a sparse solution of  $\beta$  (Zhang and Zhang, 2012), as  $\ell_1$  regularization does (Tibshirani, 1996).

### 197 3.2 Sequential LASSO and OMP

198 This connection between Sequential Attention and Sequential LASSO gives us a new perspective  
 199 about how Sequential Attention works. The only known guarantee for Sequential LASSO, to the best  
 200 of our knowledge, is a statistical recovery result when the input is a sparse linear combination with  
 201 Gaussian noise in the ultra high-dimensional setting (Luo and Chen, 2014). This does not, however,  
 202 fully explain why Sequential Attention is such an effective feature selection algorithm.

203 To bridge our main results, we prove a novel equivalence between Sequential LASSO and OMP.

204 **Theorem 3.3.** Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  be a design matrix with  $\ell_2$  unit vector columns, and let  $\mathbf{y} \in \mathbb{R}^d$  denote  
 205 the response, also an  $\ell_2$  unit vector. The Sequential LASSO algorithm maintains a set of features  
 206  $S \subseteq [d]$  such that, at each feature selection step, it selects a feature  $i \in \bar{S}$  such that

$$|\langle \mathbf{X}_i, \mathbf{P}_{\bar{S}}^\perp \mathbf{y} \rangle| = \|\mathbf{X}^\top \mathbf{P}_{\bar{S}}^\perp \mathbf{y}\|_\infty,$$

207 where  $\mathbf{X}_S$  is the  $n \times |S|$  matrix given formed by the columns of  $\mathbf{X}$  indexed by  $S$ , and  $\mathbf{P}_{\bar{S}}^\perp$  is the  
 208 projection matrix onto the orthogonal complement of the span of  $\mathbf{X}_S$ .

209 Note that this is extremely close to saying that Sequential LASSO and OMP select the exact same set  
 210 of features. The only difference appears when there are multiple features with norm  $\|\mathbf{X}^\top \mathbf{P}_{\bar{S}}^\perp \mathbf{y}\|_\infty$ .  
 211 In this case, it is possible that Sequential LASSO chooses the next feature from a set of features  
 212 that is strictly smaller than the set of features from which OMP chooses, so the ‘‘tie-breaking’’ can  
 213 differ between the two algorithms. In practice, however, this rarely happens. For instance, if only one  
 214 feature is selected at each step, which is the case with probability 1 if random continuous noise is  
 215 added to the data, then Sequential LASSO and OMP will select the exact same set of features.

216 **Remark 3.4.** It was shown in (Luo and Chen, 2014) that Sequential LASSO is equivalent to OMP in  
 217 the statistical recovery regime, i.e., when  $\mathbf{y} = \mathbf{X}\beta^* + \varepsilon$  for some true sparse weight vector  $\beta^*$  and  
 218 i.i.d. Gaussian noise  $\varepsilon \sim \mathcal{N}(0, \sigma \mathbf{I}_n)$ , under an ultra high-dimensional regime where the dimension  $d$   
 219 is exponential in the number of examples  $n$ . We prove this equivalence in the *fully general setting*.

220 The argument below shows that Sequential LASSO and OMP are equivalent, thus establishing  
 221 that regularized linear Sequential Attention and Sequential LASSO have the same approximation  
 222 guarantees as OMP.

223 **Geometry of Sequential LASSO.** We first study the geometry of optimal solutions to Equation (4).  
 224 Let  $S \subseteq [d]$  be the set of currently selected features. Following work on the LASSO in Tibshirani  
 225 and Taylor (2011), we rewrite (4) as the following constrained optimization problem:

$$\begin{aligned} \min_{\mathbf{z} \in \mathbb{R}^n, \beta \in \mathbb{R}^d} \quad & \frac{1}{2} \|\mathbf{z} - \mathbf{y}\|_2^2 + \lambda \|\beta_{\bar{S}}\|_1 \\ \text{subject to} \quad & \mathbf{z} = \mathbf{X}\beta. \end{aligned} \tag{8}$$

226 It can then be shown that the dual problem is equivalent to finding the projection, i.e., closest point in  
 227 Euclidean distance,  $\mathbf{u} \in \mathbb{R}^n$  of  $\mathbf{P}_{\bar{S}}^\perp \mathbf{y}$  onto the polyhedral section  $C_\lambda \cap \text{colspan}(\mathbf{X}_S)^\perp$ , where

$$C_\lambda := \{\mathbf{u}' \in \mathbb{R}^n : \|\mathbf{X}^\top \mathbf{u}'\|_\infty \leq \lambda\}$$

228 and  $\text{colspan}(\mathbf{X}_S)^\perp$  denotes the orthogonal complement of  $\text{colspan}(\mathbf{X}_S)$ . See Appendix A.1 for the  
 229 full details. The primal and dual variables are related through  $\mathbf{z}$  by

$$\mathbf{X}\beta = \mathbf{z} = \mathbf{y} - \mathbf{u}. \tag{9}$$

230 **Selection of features in Sequential LASSO.** Next, we analyze how Sequential LASSO selects  
 231 its features. Let  $\beta_S^* = \mathbf{X}_S^+ \mathbf{y}$  be the optimal solution for features restricted in  $S$ . Then, subtracting  
 232  $\mathbf{X}_S \beta_S^*$  from both sides of (9) gives

$$\begin{aligned} \mathbf{X}\beta - \mathbf{X}_S \beta_S^* &= \mathbf{y} - \mathbf{X}_S \beta_S^* - \mathbf{u} \\ &= \mathbf{P}_{\bar{S}}^\perp \mathbf{y} - \mathbf{u}. \end{aligned} \tag{10}$$

233 Note that if  $\lambda \geq \|\mathbf{X}^\top \mathbf{P}_{\bar{S}}^\perp \mathbf{y}\|_\infty$ , then the projection of  $\mathbf{P}_{\bar{S}}^\perp \mathbf{y}$  onto  $C_\lambda$  is just  $\mathbf{u} = \mathbf{P}_{\bar{S}}^\perp \mathbf{y}$ , so by (10),

$$\mathbf{X}\beta - \mathbf{X}_S \beta_S^* = \mathbf{P}_{\bar{S}}^\perp \mathbf{y} - \mathbf{P}_{\bar{S}}^\perp \mathbf{y} = \mathbf{0},$$

234 meaning that  $\beta$  is zero outside of  $S$ . We now show that for  $\lambda$  slightly smaller than  $\|\mathbf{X}^\top \mathbf{P}_{\bar{S}}^\perp \mathbf{y}\|_\infty$ , the  
 235 residual  $\mathbf{P}_{\bar{S}}^\perp \mathbf{y} - \mathbf{u}$  is in the span of features  $\mathbf{X}_i$  that maximize the correlation with  $\mathbf{P}_{\bar{S}}^\perp \mathbf{y}$ .

236 **Lemma 3.5** (Projection residuals of the Sequential LASSO). Let  $\mathbf{p}$  denote the projection of  $\mathbf{P}_{\bar{S}}^\perp \mathbf{y}$   
 237 onto  $C_\lambda \cap \text{colspan}(\mathbf{X}_S)^\perp$ . There exists  $\lambda_0 < \|\mathbf{X}^\top \mathbf{P}_{\bar{S}}^\perp \mathbf{y}\|_\infty$  such that for all  $\lambda \in (\lambda_0, \|\mathbf{X}^\top \mathbf{P}_{\bar{S}}^\perp \mathbf{y}\|_\infty)$   
 238 the residual  $\mathbf{P}_{\bar{S}}^\perp \mathbf{y} - \mathbf{p}$  lies on  $\text{colspan}(\mathbf{X}_T)$ , for

$$T := \{i \in [d] : |\langle \mathbf{X}_i, \mathbf{P}_{\bar{S}}^\perp \mathbf{y} \rangle| = \|\mathbf{X}^\top \mathbf{P}_{\bar{S}}^\perp \mathbf{y}\|_\infty\}.$$

239 We defer the proof of Lemma 3.5 to Appendix A.2.

240 By Lemma 3.5 and (10), the optimal  $\beta$  when selecting the next feature has the following properties:

- 241 1. if  $i \in S$ , then  $\beta_i$  is equal to the  $i$ -th value in the previous solution  $\beta_S^*$ ; and
- 242 2. if  $i \notin S$ , then  $\beta_i$  can be nonzero only if  $i \in T$ .

243 It follows that Sequential LASSO selects a feature that maximizes the correlation  $|\langle \mathbf{X}_j, \mathbf{P}_S^\perp \mathbf{y} \rangle|$ , just  
 244 as OMP does. Thus, we have shown an equivalence between Sequential LASSO and OMP without  
 245 any additional assumptions.

## 246 4 Experiments

### 247 4.1 Feature selection for neural networks

248 **Small-scale experiments.** We investigate the performance of Sequential Attention, as presented in  
 249 Algorithm 1, through experiments on standard feature selection benchmarks for neural networks. In  
 250 these experiments, we consider six datasets used in experiments in Lemhadri et al. (2021); Baln et al.  
 251 (2019), and select  $k = 50$  features using a one-layer neural network with hidden width 67 and ReLU  
 252 activation (just as in these previous works). For more points of comparison, we also implement the  
 253 attention-based feature selection algorithms of Baln et al. (2019); Liao et al. (2021) and the Group  
 254 LASSO, which has been considered in many works that aim to sparsify neural networks as discussed  
 255 in Section 1.1. We also implement natural adaptations of the Sequential LASSO and OMP for neural  
 256 networks and evaluate their performance.

257 In Figure 3, we see that Sequential Attention is competitive with or outperforms all feature selection  
 258 algorithms on this benchmark suite. For each algorithm, we report the mean of the prediction  
 259 accuracies averaged over five feature selection trials. We provide more details about the experimental  
 260 setup in Appendix B.2, including specifications about each dataset in Table 1 and the mean prediction  
 261 accuracies with their standard deviations in Table 2. We also visualize the selected features on MNIST  
 262 (i.e., pixels) in Figure 5.

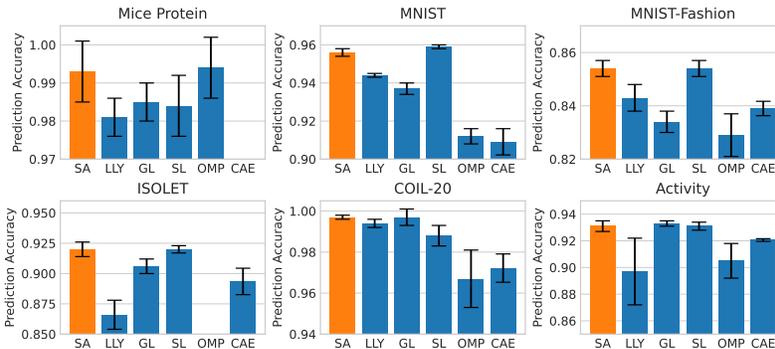


Figure 3: Feature selection results for small-scale neural network experiments. Here, SA = Sequential Attention, LLY = (Liao et al., 2021), GL = Group LASSO, SL = Sequential LASSO, OMP = OMP, and CAE = Concrete Autoencoder (Baln et al., 2019).

263 We note that our algorithm is considerably more efficient compared to prior feature selection algo-  
 264 rithms, especially those designed for neural networks. This is because many of these prior algorithms  
 265 introduce entire subnetworks to train (Baln et al., 2019; Gui et al., 2019; Wojtas and Chen, 2020; Liao  
 266 et al., 2021), whereas Sequential Attention only adds  $d$  additional trainable variables. Furthermore, in  
 267 these experiments, we implement an optimized version of Algorithm 1 that only trains one model  
 268 rather than  $k$  models, by partitioning the training epochs into  $k$  parts and selecting one feature in  
 269 each of these  $k$  parts. Combining these two aspects makes for an extremely efficient algorithm. We  
 270 provide an evaluation of the running time efficiency of Sequential Attention in Appendix B.2.3.

271 **Large-scale experiments.** To demonstrate the scalability of our algorithm, we perform large-scale  
 272 feature selection experiments on the Criteo click dataset, which consists of 39 features and over three

273 billion examples for predicting click-through rates (Diemert Eustache, Meynet Julien et al., 2017).  
 274 Our results in Figure 4 show that Sequential Attention outperforms other methods when at least 15  
 275 features are selected. In particular, these plots highlight the fact that Sequential Attention excels at  
 276 finding valuable features once a few features are already in the model, and that it has substantially less  
 277 variance than LASSO-based feature selection algorithms. See Appendix B.3 for further discussion.

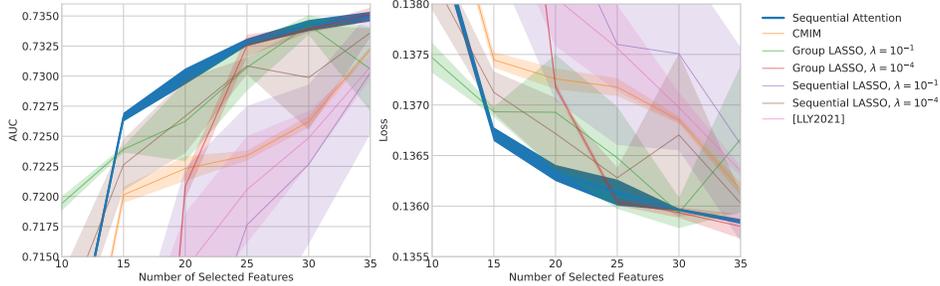


Figure 4: AUC and log loss when selecting  $k \in \{10, 15, 20, 25, 30, 35\}$  features for Criteo dataset.

## 278 4.2 The role of Hadamard product overparameterization in attention

279 In Section 1, we argued that Sequential Attention has provable guarantees for least squares linear  
 280 regression by showing that a version that removes the softmax and introduces  $\ell_2$  regularization  
 281 results in an algorithm that is equivalent to OMP. Thus, there is a gap between the implementation of  
 282 Sequential Attention in Algorithm 1 and our theoretical analysis. We empirically bridge this gap by  
 283 showing that regularized linear Sequential Attention yields results that are almost indistinguishable  
 284 to the original version. In Figure 10 (Appendix B.5), we compare the following Hadamard product  
 285 overparameterization schemes:

- 286 • softmax: as described in Section 1
- 287 •  $\ell_1$ :  $\mathbf{s}_i(\mathbf{w}) = |\mathbf{w}_i|$  for  $i \in \bar{\mathcal{S}}$ , which captures the provable variant discussed in Section 1
- 288 •  $\ell_2$ :  $\mathbf{s}_i(\mathbf{w}) = |\mathbf{w}_i|^2$  for  $i \in \bar{\mathcal{S}}$
- 289 •  $\ell_1$  normalized:  $\mathbf{s}_i(\mathbf{w}) = |\mathbf{w}_i| / \sum_{j \in \bar{\mathcal{S}}} |\mathbf{w}_j|$  for  $i \in \bar{\mathcal{S}}$
- 290 •  $\ell_2$  normalized:  $\mathbf{s}_i(\mathbf{w}) = |\mathbf{w}_i|^2 / \sum_{j \in \bar{\mathcal{S}}} |\mathbf{w}_j|^2$  for  $i \in \bar{\mathcal{S}}$

291 Further, for each of the benchmark datasets, all of these variants outperform LassoNet and the other  
 292 baselines considered in Lemhadri et al. (2021). See Appendix B.5 for more details.

## 293 5 Conclusion

294 This work introduces Sequential Attention, an adaptive attention-based feature selection algorithm  
 295 designed in part for neural networks. Empirically, Sequential Attention improves significantly upon  
 296 previous methods on widely-used benchmarks. Theoretically, we show that a relaxed variant of  
 297 Sequential Attention is equivalent to Sequential LASSO (Luo and Chen, 2014). In turn, we prove a  
 298 novel connection between Sequential LASSO and Orthogonal Matching Pursuit, thus transferring the  
 299 provable guarantees of OMP to Sequential Attention and shedding light on our empirical results. This  
 300 analysis also provides new insights into the the role of attention for feature selection via adaptivity,  
 301 overparameterization, and connections to marginal gains.

302 We conclude with a number of open questions that stem from this work. The first question concerns  
 303 the generalization of our theoretical results for Sequential LASSO to other models. OMP admits  
 304 provable guarantees for a wide class of generalized linear models (Elenberg et al., 2018), so is the  
 305 same true for Sequential LASSO? Our second question concerns the role of softmax in Algorithm 1.  
 306 Our experimental results suggest that using softmax for overparameterization may not be necessary, and  
 307 that a wide variety of alternative expressions can be used. On the other hand, our provable guarantees  
 308 only hold for the overparameterization scheme in the regularized linear Sequential Attention algorithm  
 309 (see Definition 3.1). Can we obtain a deeper understanding about the pros and cons of the softmax  
 310 and other overparameterization patterns, both theoretically and empirically?

## 311 References

- 312 Jason M. Altschuler, Aditya Bhaskara, Gang Fu, Vahab S. Mirrokni, Afshin Rostamizadeh, and  
313 Morteza Zadimoghaddam. Greedy column subset selection: New bounds and distributed algorithms.  
314 In *Proceedings of the 33rd International Conference on Machine Learning*, volume 48, pages  
315 2539–2548. JMLR, 2016.
- 316 Jose M Alvarez and Mathieu Salzmann. Learning the number of neurons in deep networks. *Advances*  
317 *in neural information processing systems*, 29, 2016.
- 318 Sercan Ö Arik and Tomas Pfister. TabNet: Attentive interpretable tabular learning. In *Proceedings of*  
319 *the AAAI Conference on Artificial Intelligence*, volume 35, pages 6679–6687, 2021.
- 320 Muhammed Fatih Balin, Abubakar Abid, and James Zou. Concrete autoencoders: Differentiable  
321 feature selection and reconstruction. In *International conference on machine learning*, pages  
322 444–453. PMLR, 2019.
- 323 Yoav Benjamini, Dan Drai, Greg Elmer, Neri Kafkafi, and Ilan Golani. Controlling the false discovery  
324 rate in behavior genetics research. *Behavioural Brain Research*, 125(1-2):279–284, 2001.
- 325 Mohamed Bennisar, Yulia Hicks, and Rossitza Setchi. Feature selection using joint mutual informa-  
326 tion maximisation. *Expert Systems with Applications*, 42(22):8520–8532, 2015.
- 327 Jeff Bilmes. Submodularity in machine learning and artificial intelligence. *arXiv preprint*  
328 *arXiv:2202.00132*, 2022.
- 329 Vadim Borisov, Johannes Haug, and Gjergji Kasneci. CancelOut: A layer for feature selection in  
330 deep neural networks. In *International conference on artificial neural networks*, pages 72–83.  
331 Springer, 2019.
- 332 Stephen Boyd and Lieven Vandenbergh. *Convex optimization*. Cambridge University Press, 2004.
- 333 Sébastien Bubeck and Mark Sellke. A universal law of robustness via isoperimetry. In *Advances in*  
334 *Neural Information Processing Systems*, pages 28811–28822, 2021.
- 335 Brais Cancela, Verónica Bolón-Canedo, Amparo Alonso-Betanzos, and João Gama. A scalable  
336 saliency-based feature selection method with instance-level information. *Knowl. Based Syst.*, 192:  
337 105326, 2020.
- 338 Emmanuel J Candes and Terence Tao. Near-optimal signal recovery from random projections:  
339 Universal encoding strategies? *IEEE Transactions on Information Theory*, 52(12):5406–5425,  
340 2006.
- 341 Jianbo Chen, Le Song, Martin Wainwright, and Michael Jordan. Learning to explain: An information-  
342 theoretic perspective on model interpretation. In *International Conference on Machine Learning*,  
343 pages 883–892. PMLR, 2018.
- 344 Lin Chen, Hossein Esfandiari, Gang Fu, Vahab S. Mirrokni, and Qian Yu. Feature Cross Search via  
345 Submodular Optimization. In *29th Annual European Symposium on Algorithms (ESA 2021)*, pages  
346 31:1–31:16, 2021.
- 347 Abhimanyu Das and David Kempe. Submodular meets spectral: Greedy algorithms for subset  
348 selection, sparse approximation and dictionary selection. In *Proceedings of the 28th International*  
349 *Conference on Machine Learning*, pages 1057–1064, 2011.
- 350 Sandipan Das, Alireza M Javid, Prakash Borpatra Gohain, Yonina C Eldar, and Saikat Chatterjee.  
351 Neural greedy pursuit for feature selection. *arXiv preprint arXiv:2207.09390*, 2022.
- 352 Diemert Eustache, Meynet Julien, Pierre Galland, and Damien Lefortier. Attribution modeling  
353 increases efficiency of bidding in display advertising. In *Proceedings of the AdKDD and TargetAd*  
354 *Workshop, KDD, Halifax, NS, Canada, August, 14, 2017*. ACM, 2017.
- 355 Chris H. Q. Ding and Hanchuan Peng. Minimum redundancy feature selection from microarray gene  
356 expression data. *J. Bioinform. Comput. Biol.*, 3(2):185–206, 2005.

- 357 David L Donoho and Michael Elad. Optimally sparse representation in general (nonorthogonal)  
358 dictionaries via  $\ell_1$  minimization. *Proceedings of the National Academy of Sciences*, 100(5):  
359 2197–2202, 2003.
- 360 Bradley Efron, Trevor Hastie, Iain Johnstone, and Robert Tibshirani. Least angle regression. *The*  
361 *Annals of Statistics*, 32(2):407–499, 2004.
- 362 Ethan R Elenberg, Rajiv Khanna, Alexandros G Dimakis, and Sahand Negahban. Restricted strong  
363 convexity implies weak submodularity. *The Annals of Statistics*, 46(6B):3539–3568, 2018.
- 364 Matthew Fahrback, Vahab Mirrokni, and Morteza Zadimoghaddam. Non-monotone submodular  
365 maximization with nearly optimal adaptivity and query complexity. In *International Conference*  
366 *on Machine Learning*, pages 1833–1842. PMLR, 2019a.
- 367 Matthew Fahrback, Vahab Mirrokni, and Morteza Zadimoghaddam. Submodular maximization with  
368 nearly optimal approximation, adaptivity and query complexity. In *Proceedings of the Thirtieth*  
369 *Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 255–273. SIAM, 2019b.
- 370 François Fleuret. Fast binary feature selection with conditional mutual information. *Journal of*  
371 *Machine learning research*, 5(9), 2004.
- 372 Ning Gui, Danni Ge, and Ziyin Hu. AFS: An attention-based mechanism for supervised feature  
373 selection. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages  
374 3705–3713, 2019.
- 375 Marwa El Halabi, Suraj Srinivas, and Simon Lacoste-Julien. Data-efficient structured pruning via  
376 submodular optimization. *arXiv preprint arXiv:2203.04940*, 2022.
- 377 Peter D Hoff. Lasso, fractional norm and structured sparse estimation using a Hadamard product  
378 parametrization. *Computational Statistics & Data Analysis*, 115:186–198, 2017.
- 379 Vadim Lebedev and Victor Lempitsky. Fast convnets using group-wise brain damage. In *Proceedings*  
380 *of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2554–2564, 2016.
- 381 Ismael Lemhadri, Feng Ruan, and Rob Tibshirani. Lassonet: Neural networks with feature sparsity.  
382 In *International Conference on Artificial Intelligence and Statistics*, pages 10–18. PMLR, 2021.
- 383 Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for  
384 efficient convnets. In *5th International Conference on Learning Representations (ICLR)*, 2017a.
- 385 Jundong Li, Kewei Cheng, Suhang Wang, Fred Morstatter, Robert P Trevino, Jiliang Tang, and Huan  
386 Liu. Feature selection: A data perspective. *ACM Computing Surveys (CSUR)*, 50(6):1–45, 2017b.
- 387 Yifeng Li, Chih-Yu Chen, and Wyeth W Wasserman. Deep feature selection: Theory and application  
388 to identify enhancers and promoters. *Journal of Computational Biology*, 23(5):322–336, 2016.
- 389 Yiwen Liao, Raphaël Latty, and Bin Yang. Feature selection using batch-wise attenuation and feature  
390 mask normalization. In *2021 International Joint Conference on Neural Networks (IJCNN)*, pages  
391 1–9. IEEE, 2021.
- 392 Christos Louizos, Max Welling, and Diederik P. Kingma. Learning sparse neural networks through  
393  $L_0$  regularization. In *6th International Conference on Learning Representations (ICLR)*, 2018.
- 394 Yang Lu, Yingying Fan, Jinchu Lv, and William Stafford Noble. DeepPINK: Reproducible feature  
395 selection in deep neural networks. *Advances in Neural Information Processing Systems*, 31, 2018.
- 396 Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in*  
397 *Neural Information Processing Systems*, 30, 2017.
- 398 Shan Luo and Zehua Chen. Sequential lasso cum EBIC for feature selection with ultra-high di-  
399 mensional feature space. *Journal of the American Statistical Association*, 109(507):1229–1240,  
400 2014.

- 401 Yagyensh Chandra Pati, Ramin Rezaifar, and Perinkulam Sambamurthy Krishnaprasad. Orthogonal  
402 matching pursuit: Recursive function approximation with applications to wavelet decomposition.  
403 In *Proceedings of 27th Asilomar conference on signals, systems and computers*, pages 40–44.  
404 IEEE, 1993.
- 405 Eric Price, Sandeep Silwal, and Samson Zhou. Hardness and algorithms for robust and sparse  
406 optimization. In *International Conference on Machine Learning*, pages 17926–17944. PMLR,  
407 2022.
- 408 Debaditya Roy, K Sri Rama Murty, and C Krishna Mohan. Feature selection using deep neural  
409 networks. In *2015 International Joint Conference on Neural Networks (IJCNN)*, pages 1–6. IEEE,  
410 2015.
- 411 Simone Scardapane, Danilo Comminiello, Amir Hussain, and Aurelio Uncini. Group sparse regular-  
412 ization for deep neural networks. *Neurocomputing*, 241:81–89, 2017.
- 413 Blaz Skrlj, Saso Dzeroski, Nada Lavrac, and Matej Petkovic. Feature importance estimation with  
414 self-attention networks. In *24th European Conference on Artificial Intelligence (ECAI)*, volume  
415 325 of *Frontiers in Artificial Intelligence and Applications*, pages 1491–1498. IOS Press, 2020.
- 416 Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. SmoothGrad:  
417 Removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- 418 Suraj Srinivas and François Fleuret. Full-gradient representation for neural network visualization.  
419 *Advances in Neural Information Processing Systems*, 32, 2019.
- 420 Suraj Srinivas, Akshayvarun Subramanya, and R Venkatesh Babu. Training sparse neural networks.  
421 In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition workshops*,  
422 pages 138–145, 2017.
- 423 Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In  
424 *International Conference on Machine Learning*, pages 3319–3328. PMLR, 2017.
- 425 Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical  
426 Society: Series B (Methodological)*, 58(1):267–288, 1996.
- 427 Ryan J Tibshirani. Equivalences between sparse models and neural networks. *Working Notes*. URL  
428 <https://www.stat.cmu.edu/~ryantibs/papers/sparsitynn.pdf>, 2021.
- 429 Ryan J. Tibshirani and Jonathan Taylor. The solution path of the generalized lasso. *The Annals of  
430 Statistics*, 39(3):1335–1371, 2011.
- 431 Andrii Trelin and Aleš Procházka. Binary stochastic filtering: Feature selection and beyond. *arXiv  
432 preprint arXiv:2007.03920*, 2020.
- 433 Tomas Vaskevicius, Varun Kanade, and Patrick Rebeschini. Implicit regularization for optimal sparse  
434 recovery. *Advances in Neural Information Processing Systems*, 32, 2019.
- 435 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
436 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing  
437 Systems*, 30, 2017.
- 438 Qian Wang, Jiaying Zhang, Sen Song, and Zheng Zhang. Attentional neural network: Feature  
439 selection using cognitive feedback. *Advances in Neural Information Processing Systems*, 27, 2014.
- 440 Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity in  
441 deep neural networks. *Advances in Neural Information Processing Systems*, 29, 2016.
- 442 Maksymilian Wojtas and Ke Chen. Feature importance ranking for deep learning. *Advances in  
443 Neural Information Processing Systems*, 33:5105–5114, 2020.
- 444 Yutaro Yamada, Ofir Lindenbaum, Sahand Negahban, and Yuval Kluger. Feature selection using  
445 stochastic gates. In *International Conference on Machine Learning*, pages 10648–10659. PMLR,  
446 2020.

- 447 Mao Ye and Yan Sun. Variable selection via penalized neural network: A drop-out-one loss approach.  
448 In *International Conference on Machine Learning*, pages 5620–5629. PMLR, 2018.
- 449 Xiang Ye, Zihang He, Heng Wang, and Yong Li. Towards understanding the effectiveness of attention  
450 mechanism. *arXiv preprint arXiv:2106.15067*, 2021.
- 451 Sepehr Abbasi Zadeh, Mehrdad Ghadiri, Vahab S. Mirrokni, and Morteza Zadimoghaddam. Scalable  
452 feature selection via distributed diversity maximization. In *Proceedings of the Thirty-First AAAI  
453 Conference on Artificial Intelligence*, pages 2876–2883. AAAI Press, 2017.
- 454 Cun-Hui Zhang and Tong Zhang. A general theory of concave regularization for high-dimensional  
455 sparse estimation problems. *Statist. Sci.*, 27(4):576–593, 2012.
- 456 Lei Zhao, Qinghua Hu, and Wenwu Wang. Heterogeneous feature selection with multi-modal deep  
457 neural networks and sparse group lasso. *IEEE Transactions on Multimedia*, 17(11):1936–1948,  
458 2015.
- 459 Peng Zhao, Yun Yang, and Qiao-Chu He. High-dimensional linear regression via implicit regulariza-  
460 tion. *Biometrika*, 2022.
- 461 Qin Zou, Lihao Ni, Tong Zhang, and Qian Wang. Deep learning based feature selection for remote  
462 sensing scene classification. *IEEE Geosci. Remote. Sens. Lett.*, 12(11):2321–2325, 2015.

463 **A Missing proofs from Section 3**

464 **A.1 Lagrangian dual of Sequential LASSO**

465 We first show that the Lagrangian dual of (8) is equivalent to the following problem:

$$\begin{aligned} \min_{\mathbf{u} \in \mathbb{R}^n} \quad & \frac{1}{2} \|\mathbf{y} - \mathbf{u}\|_2^2 \\ \text{subject to} \quad & \|\mathbf{X}^\top \mathbf{u}\|_\infty \leq \lambda \\ & \mathbf{X}_j^\top \mathbf{u} = 0, \quad \forall j \in S \end{aligned} \quad (11)$$

466 We then use the Pythagorean theorem to replace  $\mathbf{y}$  by  $\mathbf{P}_S^\perp \mathbf{y}$ .

467 First consider the Lagrangian dual problem:

$$\max_{\mathbf{u} \in \mathbb{R}^n} \min_{\mathbf{z} \in \mathbb{R}^n, \beta \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{z} - \mathbf{y}\|_2^2 + \lambda \|\beta\|_1 + \mathbf{u}^\top (\mathbf{z} - \mathbf{X}\beta). \quad (12)$$

468 Note that the primal problem is strictly feasible and convex, so strong duality holds (see, e.g., Section  
469 5.2.3 of [Boyd and Vandenberghe \(2004\)](#)). Considering just the terms involving the variable  $\mathbf{z}$  in (12),  
470 we have that

$$\begin{aligned} \frac{1}{2} \|\mathbf{z} - \mathbf{y}\|_2^2 + \mathbf{u}^\top \mathbf{z} &= \frac{1}{2} \|\mathbf{z}\|_2^2 - (\mathbf{y} - \mathbf{u})^\top \mathbf{z} + \frac{1}{2} \|\mathbf{y}\|_2^2 \\ &= \frac{1}{2} \|\mathbf{z} - (\mathbf{y} - \mathbf{u})\|_2^2 + \frac{1}{2} \|\mathbf{y}\|_2^2 - \frac{1}{2} \|\mathbf{y} - \mathbf{u}\|_2^2, \end{aligned}$$

471 which is minimized at  $\mathbf{z} = \mathbf{y} - \mathbf{u}$  as  $\mathbf{z}$  varies over  $\mathbb{R}^n$ . On the other hand, consider just the terms  
472 involving the variable  $\beta$  in (12), that is,

$$\lambda \|\beta\|_1 - \mathbf{u}^\top \mathbf{X}\beta. \quad (13)$$

473 Note that if  $\mathbf{X}^\top \mathbf{u}$  is nonzero on any coordinate in  $S$ , then (13) can be made arbitrarily negative  
474 by setting  $\beta_S$  to be zero and  $\beta_{S^c}$  appropriately. Similarly, if  $\|\mathbf{X}^\top \mathbf{u}\|_\infty > \lambda$ , then (13) can also be  
475 made to be arbitrarily negative. On the other hand, if  $(\mathbf{X}^\top \mathbf{u})_S = \mathbf{0}$  and  $\|\mathbf{X}^\top \mathbf{u}\|_\infty \leq \lambda$ , then (13) is  
476 minimized at 0. This gives the dual in Equation (11).

477 We now show that by the Pythagorean theorem, we can project  $\mathbf{P}_S^\perp \mathbf{y}$  in (11) rather than  $\mathbf{y}$ . In (11),  
478 recall that  $\mathbf{u}$  is constrained to be in  $\text{colspan}(\mathbf{X}_S)^\perp$ . Then, by the Pythagorean theorem, we have

$$\begin{aligned} \frac{1}{2} \|\mathbf{y} - \mathbf{u}\|_2^2 &= \frac{1}{2} \|\mathbf{y} - \mathbf{P}_S^\perp \mathbf{y} + \mathbf{P}_S^\perp \mathbf{y} - \mathbf{u}\|_2^2 \\ &= \frac{1}{2} \|\mathbf{y} - \mathbf{P}_S^\perp \mathbf{y}\|_2^2 + \frac{1}{2} \|\mathbf{P}_S^\perp \mathbf{y} - \mathbf{u}\|_2^2, \end{aligned}$$

479 since  $\mathbf{y} - \mathbf{P}_S^\perp \mathbf{y} = \mathbf{P}_S \mathbf{y}$  is orthogonal to  $\text{colspan}(\mathbf{X}_S)^\perp$  and both  $\mathbf{P}_S^\perp \mathbf{y}$  and  $\mathbf{u}$  are in  $\text{colspan}(\mathbf{X}_S)^\perp$ .  
480 The first term in the above does not depend on  $\mathbf{u}$  and thus we may discard it. Our problem therefore  
481 reduces to projecting  $\mathbf{P}_S^\perp \mathbf{y}$  onto  $C_\lambda \cap \text{colspan}(\mathbf{X}_S)^\perp$ , rather than  $\mathbf{y}$ .

482 **A.2 Proof of Lemma 3.5**

483 *Proof of Lemma 3.5.* Our approach is to reduce the projection of  $\mathbf{P}_S^\perp \mathbf{y}$  onto the polytope defined by  
484  $C_\lambda \cap \text{colspan}(\mathbf{X})^\perp$  to a projection onto an affine space.

485 We first argue that it suffices to project onto the faces of  $C_\lambda$  specified by set  $T$ . For  $\lambda > 0$ , feature  
486 indices  $i \in [d]$ , and signs  $\pm$ , we define the faces

$$F_{\lambda, i, \pm} := \{\mathbf{u} \in \mathbb{R}^n : \pm \langle \mathbf{X}_i, \mathbf{u} \rangle = \lambda\}$$

487 of  $C_\lambda$ . Let  $\lambda = (1 - \varepsilon) \|\mathbf{X}^\top \mathbf{P}_S^\perp \mathbf{y}\|_\infty$ , for  $\varepsilon > 0$  to be chosen sufficiently small. Then clearly

$$(1 - \varepsilon) \mathbf{P}_S^\perp \mathbf{y} \in C_\lambda \cap \text{colspan}(\mathbf{X}_S)^\perp,$$

488 so

$$\min_{\mathbf{u} \in C_\lambda \cap \text{colspan}(\mathbf{X}_S)^\perp} \|\mathbf{P}_S^\perp \mathbf{y} - \mathbf{u}\|_2^2 \leq \|\mathbf{P}_S^\perp \mathbf{y} - (1 - \varepsilon) \mathbf{P}_S^\perp \mathbf{y}\|_2^2$$

$$= \varepsilon^2 \|\mathbf{P}_S^\perp \mathbf{y}\|_2^2.$$

489 In fact,  $(1 - \varepsilon)\mathbf{P}_S^\perp \mathbf{y}$  lies on the intersection of faces  $F_{\lambda, i, \pm}$  for an appropriate choice of signs and  
 490  $i \in T$ . Without loss of generality, we assume that these faces are just  $F_{\lambda, i, +}$  for  $i \in T$ . Note also that  
 491 for any  $i \notin T$ ,

$$\begin{aligned} \min_{\mathbf{u} \in F_{\lambda, i, \pm}} \|\mathbf{P}_S^\perp \mathbf{y} - \mathbf{u}\|_2^2 &\geq \min_{\mathbf{u} \in F_{\lambda, i, \pm}} \langle \mathbf{X}_i, \mathbf{P}_S^\perp \mathbf{y} - \mathbf{u} \rangle^2 && \text{(Cauchy-Schwarz, } \|\mathbf{X}_i\|_2 \leq 1) \\ &= \min_{\mathbf{u} \in F_{\lambda, i, \pm}} |\mathbf{X}_i^\top \mathbf{P}_S^\perp \mathbf{y} - \mathbf{X}_i^\top \mathbf{u}|^2 \\ &= (|\mathbf{X}_i^\top \mathbf{P}_S^\perp \mathbf{y}| - \lambda)^2 && (\mathbf{u} \in F_{\lambda, i, \pm}) \\ &\geq \left( (1 - \varepsilon) \|\mathbf{X}^\top \mathbf{P}_S^\perp \mathbf{y}\|_\infty - \|\mathbf{X}_T^\top \mathbf{P}_S^\perp \mathbf{y}\|_\infty \right)^2. \end{aligned}$$

492 For all  $\varepsilon < \varepsilon_0$ , for  $\varepsilon_0$  small enough, this is larger than  $\varepsilon^2 \|\mathbf{P}_S^\perp \mathbf{y}\|_2^2$ . Thus, for  $\varepsilon$  small enough,  $\mathbf{P}_S^\perp \mathbf{y}$  is  
 493 closer to the faces  $F_{\lambda, i, +}$  for  $i \in T$  than any other face. Therefore, we set  $\lambda_0 = (1 - \varepsilon_0) \|\mathbf{X}^\top \mathbf{P}_S^\perp \mathbf{y}\|_\infty$ .

494 Now, by the complementary slackness of the KKT conditions for the projection  $\mathbf{u}$  of  $\mathbf{P}_S^\perp \mathbf{y}$  onto  $C_\lambda$ ,  
 495 for each face of  $C_\lambda$  we either have that  $\mathbf{u}$  lies on the face or that the projection does not change if we  
 496 remove the face. For  $i \notin T$ , note that by the above calculation, the projection  $\mathbf{u}$  cannot lie on  $F_{\lambda, i, \pm}$ ,  
 497 so  $\mathbf{u}$  is simply the projection onto

$$C' = \{\mathbf{u} \in \mathbb{R}^n : \mathbf{X}_T^\top \mathbf{u} \leq \lambda \mathbf{1}_T\}.$$

498 By reversing the dual problem reasoning from before, the residual of the projection onto  $C'$  must lie  
 499 on the column span of  $\mathbf{X}_T$ .  $\square$

### 500 A.3 Parameterization patterns and regularization

501 *Proof of Lemma 3.2.* The optimization problem on the left-hand side of Equation (7) with respect  
 502 to  $\beta$  is equivalent to

$$\inf_{\beta \in \mathbb{R}^d} \left( \|\mathbf{X}\beta - \mathbf{y}\|_2^2 + \frac{\lambda}{2} \inf_{\mathbf{w} \in \mathbb{R}^d} \left( \|\mathbf{w}\|_2^2 + \sum_{i \in \bar{S}} \frac{\beta_i^2}{\mathbf{s}_i(\mathbf{w})^2} \right) \right). \quad (14)$$

503 If we define

$$\tilde{Q}^*(\mathbf{x}) = \inf_{\mathbf{w} \in \mathbb{R}^d} \left( \|\mathbf{w}\|_2^2 + \sum_{i \in \bar{S}} \frac{\mathbf{x}_i}{\mathbf{s}_i(\mathbf{w})^2} \right),$$

504 then the LHS of (7) and (14) are equivalent to  $\inf_{\beta \in \mathbb{R}^d} (\|\mathbf{X}\beta - \mathbf{y}\|_2^2 + \frac{\lambda}{2} \tilde{Q}^*(\beta \circ \beta))$ . Re-parameterizing  
 505 the minimization problem in the definition of  $\tilde{Q}^*(\mathbf{x})$  (by setting  $\mathbf{q} = D(\mathbf{w})$ ), we obtain  $\tilde{Q}^* = Q^*$ .  $\square$

## 506 B Additional experiments

### 507 B.1 Visualization of selected MNIST features

508 In Figure 5, we present visualizations of the features (i.e., pixels) selected by Sequential Attention  
 509 and the baseline algorithms. This provides some intuition on the nature of the features that these  
 510 algorithms select. Similar visualizations for MNIST can be found in works such as [Balm et al. \(2019\)](#);  
 511 [Gui et al. \(2019\)](#); [Wojtas and Chen \(2020\)](#); [Lemhadri et al. \(2021\)](#); [Liao et al. \(2021\)](#). Note that these  
 512 visualizations serve as a basic sanity check about the kinds of pixels that these algorithms select. For  
 513 instance, the degree to which the selected pixels are “clustered” can be used to informally assess  
 514 the redundancy of features selected for image datasets, since neighboring pixels tend to represent  
 515 redundant information. It is also useful at time to assess which regions of the image are selected. For  
 516 example, the central regions of the MNIST images are more informative than the edges.

517 Sequential Attention selects a highly diverse set of pixels due to its adaptivity. Sequential LASSO also  
 518 selects a very similar set of pixels, as suggested by our theoretical analysis in Section 3. Curiously,

519 OMP does not yield a competitive set of pixels, which demonstrates that OMP does not generalize  
 520 well from least squares regression and generalized linear models to deep neural networks.

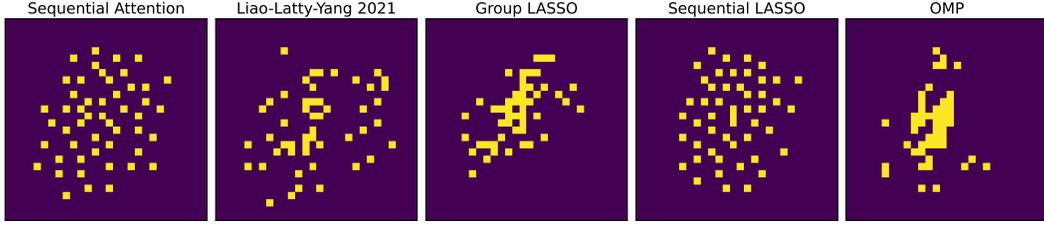


Figure 5: Visualizations of the  $k = 50$  pixels selected by the feature selection algorithms on MNIST.

## 521 B.2 Additional details on small-scale experiments

522 We start by presenting details about each of the datasets used for neural network feature selection  
 523 in [Bain et al. \(2019\)](#) and [Lemhadri et al. \(2021\)](#) in Table 1.

Table 1: Statistics about benchmark datasets.

Dataset	# Examples	# Features	# Classes	Type
Mice Protein	1,080	77	8	Biology
MNIST	60,000	784	10	Image
MNIST-Fashion	60,000	784	10	Image
ISOLET	7,797	617	26	Speech
COIL-20	1,440	400	20	Image
Activity	5,744	561	6	Sensor

524 In Figure 3, the error bars are computed using the standard deviation over five runs of the algorithm  
 525 with different random seeds. The values used to generate these plots are provided below in Table 2.

Table 2: Feature selection results for small-scale datasets (see Figure 3 for a key). These values are the average prediction accuracies on the test data and their standard deviations.

Dataset	SA	LLY	GL	SL	OMP	CAE
Mice Protein	$0.993 \pm 0.008$	$0.981 \pm 0.005$	$0.985 \pm 0.005$	$0.984 \pm 0.008$	$0.994 \pm 0.008$	$0.956 \pm 0.012$
MNIST	$0.956 \pm 0.002$	$0.944 \pm 0.001$	$0.937 \pm 0.003$	$0.959 \pm 0.001$	$0.912 \pm 0.004$	$0.909 \pm 0.007$
MNIST-Fashion	$0.854 \pm 0.003$	$0.843 \pm 0.005$	$0.834 \pm 0.004$	$0.854 \pm 0.003$	$0.829 \pm 0.008$	$0.839 \pm 0.003$
ISOLET	$0.920 \pm 0.006$	$0.866 \pm 0.012$	$0.906 \pm 0.006$	$0.920 \pm 0.003$	$0.727 \pm 0.026$	$0.893 \pm 0.011$
COIL-20	$0.997 \pm 0.001$	$0.994 \pm 0.002$	$0.997 \pm 0.004$	$0.988 \pm 0.005$	$0.967 \pm 0.014$	$0.972 \pm 0.007$
Activity	$0.931 \pm 0.004$	$0.897 \pm 0.025$	$0.933 \pm 0.002$	$0.931 \pm 0.003$	$0.905 \pm 0.013$	$0.921 \pm 0.001$

### 526 B.2.1 Model accuracies with all features

527 To adjust for the differences between the values reported in [Lemhadri et al. \(2021\)](#) and ours due (e.g.,  
 528 due to factors such as the implementation framework), we list the accuracies obtained by training the  
 529 models with all of the available features in Table 3.

### 530 B.2.2 Generalizing OMP to neural networks

531 As stated in Algorithm 2, it may be difficult to see exactly how OMP generalizes from a linear  
 532 regression model to neural networks. To do this, first observe that OMP naturally generalizes to  
 533 generalized linear models (GLMs) via the gradient of the link function, as shown in [Elenberg et al.  
 534 \(2018\)](#). Then, to extend this to neural networks, we view the neural network as a GLM for any fixing  
 535 of the hidden layer weights, and then we use the gradient of this GLM with respect to the inputs as  
 536 the feature importance scores.

Table 3: Model accuracies when trained using all available features.

Dataset	Lemhadri et al. (2021)	This paper
Mice Protein	0.990	0.963
MNIST	0.928	0.953
MNIST-Fashion	0.833	0.869
ISOLET	0.953	0.961
COIL-20	0.996	0.986
Activity	0.853	0.954

537 **B.2.3 Efficiency evaluation**

538 In this subsection, we evaluate the efficiency of the Sequential Attention algorithm against our other  
 539 baseline algorithms. We do so by fixing the number of epochs and batch size for all of the algorithms,  
 540 and then evaluating the accuracy as well as the wall clock time of each algorithm. Figures 6 and 7  
 541 provide a visualization of the accuracy and wall clock time of feature selection, while Tables 5 and 6  
 542 provide the average and standard deviations. Table 4 provides the epochs and batch size settings that  
 543 were fixed for these experiments.

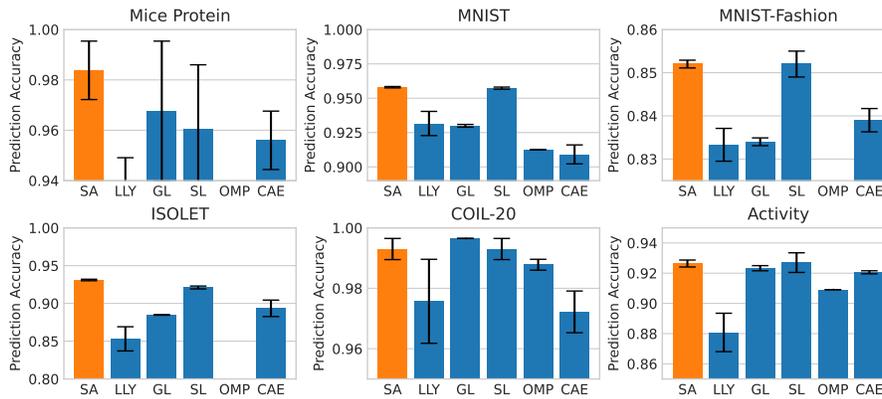


Figure 6: Feature selection accuracy for efficiency evaluation.

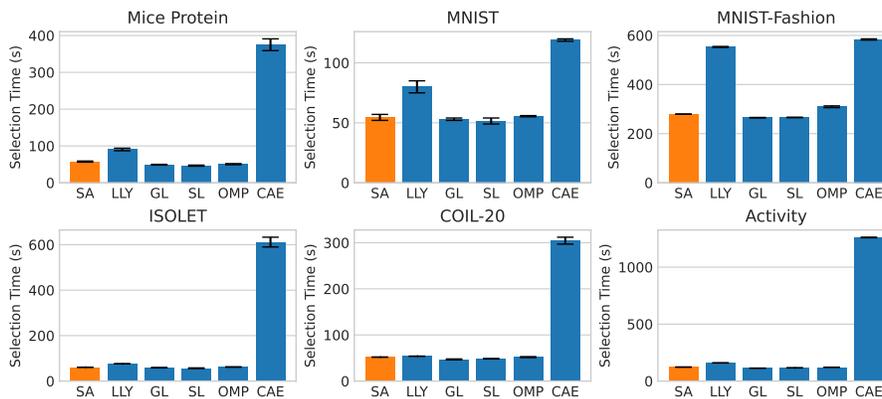


Figure 7: Feature selection wall clock time in seconds for efficiency evaluation.

544 **B.2.4 Notes on the one-pass implementation**

545 We make several remarks about the one-pass implementation of Sequential Attention. First, as noted  
 546 in Section 4.1, our practical implementation of Sequential Attention only trains one model instead

Table 4: Epochs and batch size used to compare the efficiency of feature selection algorithms.

Dataset	Epochs	Batch Size
Mice Protein	2000	256
MNIST	50	256
MNIST-Fashion	250	128
ISOLET	500	256
COIL-20	1000	256
Activity	1000	512

Table 5: Feature selection accuracy for efficiency evaluation. We report the mean accuracy on the test dataset and the the standard deviation across five trials.

Dataset	SA	LLY	GL	SL	OMP	CAE
Mice Protein	<b>0.984 ± 0.012</b>	0.907 ± 0.042	0.968 ± 0.028	0.961 ± 0.025	0.556 ± 0.032	0.956 ± 0.012
MNIST	<b>0.958 ± 0.001</b>	0.932 ± 0.009	0.930 ± 0.001	<b>0.957 ± 0.001</b>	0.912 ± 0.000	0.909 ± 0.007
MNIST-Fashion	<b>0.852 ± 0.001</b>	0.833 ± 0.004	0.834 ± 0.001	<b>0.852 ± 0.003</b>	0.722 ± 0.029	0.839 ± 0.003
ISOLET	<b>0.931 ± 0.001</b>	0.853 ± 0.016	0.885 ± 0.000	0.921 ± 0.002	0.580 ± 0.025	0.893 ± 0.011
COIL-20	<b>0.993 ± 0.004</b>	0.976 ± 0.014	0.997 ± 0.000	<b>0.993 ± 0.004</b>	0.988 ± 0.002	0.972 ± 0.007
Activity	<b>0.926 ± 0.002</b>	0.881 ± 0.013	0.923 ± 0.002	<b>0.927 ± 0.006</b>	0.909 ± 0.000	0.921 ± 0.001

of  $k$  models. We do this by partitioning the training epochs into  $k$  parts and selecting one part in each phase. This clearly gives a more efficient running time than training  $k$  separate models. Similarly, we allow for a “warm-up” period prior to the feature selection phase, in which a small fraction of the training epochs are allotted to training just the neural network weights. When we do this one-pass implementation, we observe that it is important to reset the attention weights after each of the sequential feature selection phases, but resetting the neural network weights is not crucial for good performance.

Second, we note that when there is a large number of candidate features  $d$ , the softmax mask severely scales down the gradient updates to the model weights, which can lead to inefficient training. In these cases, it becomes important to prevent this by either using a temperature parameter in the softmax to counteract the small softmax values or by adjusting the learning rate to be high enough. Note that these two approaches can be considered to be equivalent.

### B.3 Large-scale experiments

In this section, we provide more details and discussion on our Criteo large dataset results. For this experiment, we use a dense neural network with 768, 256, and 128 neurons in each of the three hidden layers with ReLU activations. In Figure 4, the error bars are generated as the standard deviation over running the algorithm three times with different random seeds, and the shadowed regions linearly interpolate between these error bars. The values used to generate the plot are provided in Table 7 and Table 8.

We first note that this dataset is so large that it is expensive to make multiple passes through the dataset. Therefore, we modify the algorithms (both Sequential Attention and the other baselines) to make only one pass through the data by using disjoint fractions of the data for different “steps” of the algorithm. Hence, we select  $k$  features while only “training” one model.

### B.4 The role of adaptivity

We show in this section the effect of varying adaptivity on the quality of selected features in Sequential Attention. In the following experiments, we select 64 features on six datasets by selecting  $2^i$  features at a time over a fixed number of epochs of training, for  $i \in \{0, 1, 2, 3, 4, 5, 6\}$ . That is, we investigate the following question: for a fixed budget of training epochs, what is the best way to allocate the training epochs over the rounds of the feature selection process? For most datasets, we find that feature selection quality decreases as we select more features at once. An exception is the mice protein dataset, which exhibits the opposite trend, perhaps indicating that the features in the mice protein dataset are less redundant than in other datasets. Our results are summarized in Table 8 and Table 9. We also illustrate the effect of adaptivity for Sequential Attention on MNIST in Figure 9.

Table 6: Feature selection wall clock time in seconds for efficiency evaluations. These values are the mean wall clock time on the test dataset and their standard deviation across five trials.

Dataset	SA	LLY	GL	SL	OMP	CAE
Mice Protein	57.5 ± 1.5	90.5 ± 3.5	49.0 ± 1.0	46.5 ± 1.5	50.5 ± 1.5	375.0 ± 16.0
MNIST	54.5 ± 2.5	80.0 ± 5.0	53.0 ± 1.0	51.5 ± 2.5	55.5 ± 0.5	119.0 ± 1.0
MNIST-Fashion	279.5 ± 0.5	553.0 ± 2.0	265.0 ± 1.0	266.0 ± 1.0	309.5 ± 3.5	583.5 ± 2.5
ISOLET	61.0 ± 0.0	76.0 ± 0.0	59.0 ± 1.0	56.5 ± 1.5	62.0 ± 1.0	611.5 ± 21.5
COIL-20	52.0 ± 0.0	54.0 ± 0.0	47.0 ± 1.0	48.5 ± 0.5	52.0 ± 1.0	304.5 ± 7.5
Activity	123.0 ± 1.0	159.5 ± 0.5	113.5 ± 0.5	116.0 ± 0.0	121.5 ± 0.5	1260.5 ± 2.5

Table 7: AUC of Criteo large experiments. SA is Sequential Attention, GL is generalized LASSO, and SL is Sequential LASSO. The values in the header for the LASSO methods are the  $\ell_1$  regularization strengths used for each method.

$k$	SA	CMIM	GL ( $\lambda = 10^{-1}$ )	GL ( $\lambda = 10^{-4}$ )	SL ( $\lambda = 10^{-1}$ )	SL ( $\lambda = 10^{-4}$ )	Liao et al. (2021)
5	0.67232 ± 0.00015	0.63950 ± 0.00076	0.68342 ± 0.00585	0.50161 ± 0.00227	0.60278 ± 0.04473	0.67710 ± 0.00873	0.58300 ± 0.06360
10	0.70167 ± 0.00060	0.69402 ± 0.00052	0.71942 ± 0.00059	0.64262 ± 0.00187	0.62263 ± 0.06097	0.70964 ± 0.00385	0.68103 ± 0.00137
15	0.72659 ± 0.00036	0.72014 ± 0.00067	0.72392 ± 0.00027	0.65977 ± 0.00125	0.66203 ± 0.04319	0.72264 ± 0.00213	0.69762 ± 0.00654
20	0.72997 ± 0.00066	0.72232 ± 0.00103	0.72624 ± 0.00330	0.72085 ± 0.00106	0.70252 ± 0.01985	0.72668 ± 0.00307	0.71395 ± 0.00467
25	0.73281 ± 0.00030	0.72339 ± 0.00042	0.73072 ± 0.00193	0.73253 ± 0.00091	0.71764 ± 0.00987	0.73084 ± 0.00070	0.72057 ± 0.00444
30	0.73420 ± 0.00046	0.72622 ± 0.00049	0.73425 ± 0.00081	0.73390 ± 0.00026	0.72267 ± 0.00663	0.72988 ± 0.00434	0.72487 ± 0.00223
35	0.73495 ± 0.00040	0.73225 ± 0.00024	0.73058 ± 0.00350	0.73512 ± 0.00058	0.73029 ± 0.00509	0.73361 ± 0.00037	0.73078 ± 0.00102

580 One observes that the selected pixels “clump together” as  $i$  increases, indicating a greater degree of  
 581 redundancy.

582 Our empirical results in this section suggest that adaptivity greatly enhances the quality of features  
 583 selected by Sequential Attention, and in feature selection algorithms more broadly.

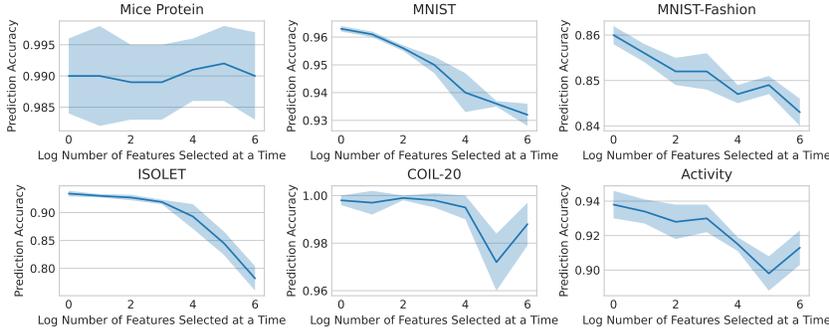


Figure 8: Sequential Attention with varying levels of adaptivity. We select 64 features for each model, and take  $2^i$  features in each round for increasing values of  $i$ . We plot accuracy as a function of  $i$ .

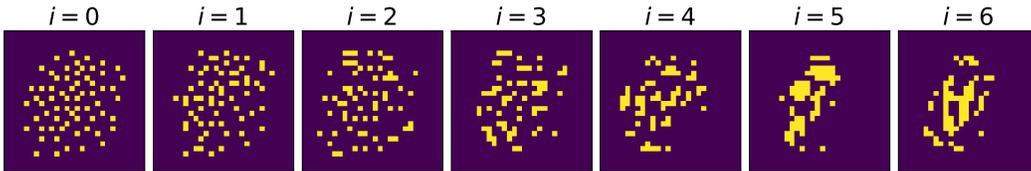


Figure 9: Sequential Attention with varying levels of adaptivity on the MNIST dataset. We select 64 features for each model, and select  $2^i$  features in each round for increasing values of  $i$ .

584 **B.5 Variations on Hadamard product parameterization**

585 We provide evaluations for different variations of the Hadamard product parameterization pattern as  
 586 described in Section 4.2. In Table 10, we provide the numerical values of the accuracies achieved.

Table 8: Log-loss of Criteo experiments. SA is Sequential Attention, GL is generalized LASSO, and SL is Sequential LASSO. The values in the header for the LASSO methods are the  $\ell_1$  regularization strengths used for each method.

$k$	SA	CMIM	GL ( $\lambda = 10^{-1}$ )	GL ( $\lambda = 10^{-4}$ )	SL ( $\lambda = 10^{-1}$ )	SL ( $\lambda = 10^{-4}$ )	Liao et al. (2021)
5	0.14123 $\pm$ 0.00005	0.14323 $\pm$ 0.00010	0.14036 $\pm$ 0.00046	0.14519 $\pm$ 0.00000	0.14375 $\pm$ 0.00163	0.14073 $\pm$ 0.00061	0.14415 $\pm$ 0.00146
10	0.13883 $\pm$ 0.00009	0.13965 $\pm$ 0.00008	0.13747 $\pm$ 0.00015	0.14339 $\pm$ 0.00019	0.14263 $\pm$ 0.00304	0.13826 $\pm$ 0.00032	0.14082 $\pm$ 0.00011
15	0.13671 $\pm$ 0.00007	0.13745 $\pm$ 0.00008	0.13693 $\pm$ 0.00005	0.14227 $\pm$ 0.00021	0.14166 $\pm$ 0.00322	0.13713 $\pm$ 0.00021	0.13947 $\pm$ 0.00050
20	0.13633 $\pm$ 0.00008	0.13726 $\pm$ 0.00010	0.13693 $\pm$ 0.00057	0.13718 $\pm$ 0.00004	0.13891 $\pm$ 0.00187	0.13672 $\pm$ 0.00035	0.13806 $\pm$ 0.00048
25	0.13613 $\pm$ 0.00013	0.13718 $\pm$ 0.00009	0.13648 $\pm$ 0.00051	0.13604 $\pm$ 0.00004	0.13760 $\pm$ 0.00099	0.13628 $\pm$ 0.00010	0.13756 $\pm$ 0.00043
30	0.13596 $\pm$ 0.00001	0.13685 $\pm$ 0.00004	0.13593 $\pm$ 0.00015	0.13594 $\pm$ 0.00005	0.13751 $\pm$ 0.00095	0.13670 $\pm$ 0.00080	0.13697 $\pm$ 0.00015
35	0.13585 $\pm$ 0.00002	0.13617 $\pm$ 0.00006	0.13666 $\pm$ 0.00073	0.13580 $\pm$ 0.00012	0.13661 $\pm$ 0.00096	0.13603 $\pm$ 0.00010	0.13635 $\pm$ 0.00005

Table 9: Sequential Attention with varying levels of adaptivity. We select 64 features for each model, and take  $2^i$  features in each round for increasing values of  $i$ . We give the accuracy as a function of  $i$ .

Dataset	$i = 0$	$i = 1$	$i = 2$	$i = 3$	$i = 4$	$i = 5$	$i = 6$
Mice Protein	0.990 $\pm$ 0.006	0.990 $\pm$ 0.008	0.989 $\pm$ 0.006	0.989 $\pm$ 0.006	0.991 $\pm$ 0.005	0.992 $\pm$ 0.006	0.990 $\pm$ 0.007
MNIST	0.963 $\pm$ 0.001	0.961 $\pm$ 0.001	0.956 $\pm$ 0.001	0.950 $\pm$ 0.003	0.940 $\pm$ 0.007	0.936 $\pm$ 0.001	0.932 $\pm$ 0.004
MNIST-Fashion	0.860 $\pm$ 0.002	0.856 $\pm$ 0.002	0.852 $\pm$ 0.003	0.852 $\pm$ 0.004	0.847 $\pm$ 0.002	0.849 $\pm$ 0.002	0.843 $\pm$ 0.003
ISOLET	0.934 $\pm$ 0.005	0.930 $\pm$ 0.003	0.927 $\pm$ 0.005	0.919 $\pm$ 0.004	0.893 $\pm$ 0.022	0.845 $\pm$ 0.021	0.782 $\pm$ 0.022
COIL-20	0.998 $\pm$ 0.002	0.997 $\pm$ 0.005	0.999 $\pm$ 0.001	0.998 $\pm$ 0.003	0.995 $\pm$ 0.005	0.972 $\pm$ 0.012	0.988 $\pm$ 0.009
Activity	0.938 $\pm$ 0.008	0.934 $\pm$ 0.007	0.928 $\pm$ 0.010	0.930 $\pm$ 0.008	0.915 $\pm$ 0.004	0.898 $\pm$ 0.010	0.913 $\pm$ 0.010

Table 10: Accuracies of Sequential Attention for different Hadamard product parameterizations.

Dataset	Softmax	$\ell_1$	$\ell_2$	$\ell_1$ normalized	$\ell_2$ normalized
Mice Protein	0.990 $\pm$ 0.006	0.993 $\pm$ 0.010	0.993 $\pm$ 0.010	0.994 $\pm$ 0.006	0.988 $\pm$ 0.008
MNIST	0.958 $\pm$ 0.002	0.957 $\pm$ 0.001	0.958 $\pm$ 0.002	0.958 $\pm$ 0.001	0.957 $\pm$ 0.001
MNIST-Fashion	0.850 $\pm$ 0.002	0.843 $\pm$ 0.004	0.850 $\pm$ 0.003	0.853 $\pm$ 0.001	0.852 $\pm$ 0.002
ISOLET	0.920 $\pm$ 0.003	0.894 $\pm$ 0.014	0.908 $\pm$ 0.009	0.921 $\pm$ 0.003	0.921 $\pm$ 0.003
COIL-20	0.997 $\pm$ 0.004	0.997 $\pm$ 0.004	0.995 $\pm$ 0.006	0.996 $\pm$ 0.005	0.996 $\pm$ 0.004
Activity	0.922 $\pm$ 0.005	0.906 $\pm$ 0.015	0.908 $\pm$ 0.012	0.933 $\pm$ 0.010	0.935 $\pm$ 0.007

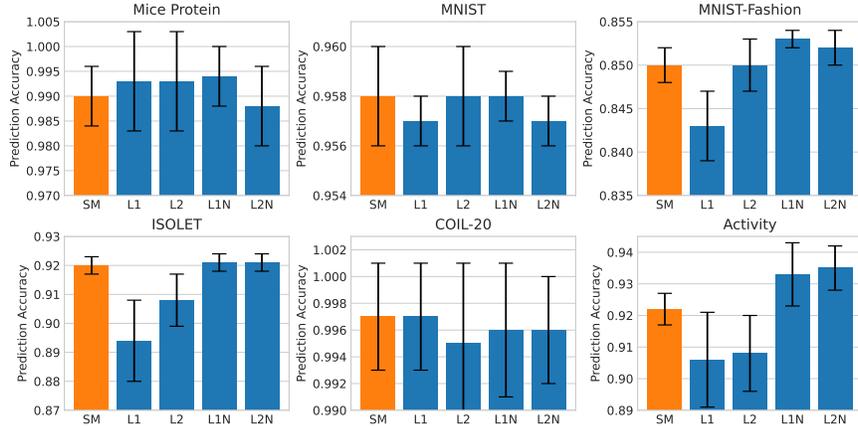


Figure 10: Accuracies of Sequential Attention for different Hadamard product parameterization patterns. Here, SM = softmax, L1 =  $\ell_1$ , L2 =  $\ell_2$ , L1N =  $\ell_1$  normalized, and L2N =  $\ell_2$  normalized.

## 587 B.6 Approximation of marginal gains

588 Finally, we present our experimental results that show the correlations between weights computed by  
 589 Sequential Attention and traditional feature selection *marginal gains*.

590 **Definition B.1** (Marginal gains). Let  $\ell : 2^{[d]} \rightarrow \mathbb{R}$  be a loss function defined on the ground set  $[d]$ .  
 591 Then, for a set  $S \subseteq [n]$  and  $i \notin S$ , the *marginal gain of  $i$  with respect to  $S$*  is  $\ell(S) - \ell(S \cup \{i\})$ .

592 In the setting of feature selection, marginal gains are often considered for measuring the importance of  
 593 candidate features  $i$  given a set  $S$  of features that have already been selected by using the set function  $\ell$ ,

594 which corresponds to the model loss when trained on a subset of features. It is known that greedily  
 595 selecting features based on their marginal gains performs well in both theory (Das and Kempe, 2011;  
 596 Elenberg et al., 2018) and practice (Das et al., 2022). These scores, however, can be extremely  
 597 expensive to compute since they require training a model for every feature considered.

598 In this experiment, we first compute the top  $k$  features selected by Sequential Attention for  $k \in$   
 599  $\{0, 9, 49\}$  on the MNIST dataset. Then we compute (1) the true marginal gains and (2) the attention  
 600 weights according to Sequential Attention, conditioned on these features being in the model. The  
 601 Sequential Attention weights are computed by only applying the attention softmax mask over the  
 602  $d - k$  unselected features, while the marginal gains are computed by explicitly training a model for  
 603 each candidate feature to be added to the preselected  $k$  features. Because our Sequential Attention  
 604 algorithm is motivated by an efficient implementation of the greedy selection algorithm that uses  
 605 marginal gains (see Section 1), one might expect that these two sets of scores are correlated in some  
 606 sense. We show this by plotting the top scores according to the two sets of scores and by computing  
 607 the Spearman correlation coefficient between the marginal gains and attention logits.

608 In the first and second rows of Figure 11, we see that the top 50 pixels according to the marginal gains  
 609 and attention weights are visually similar, avoiding previously selected regions and finding new areas  
 610 which are now important. In the third row, we quantify their similarity via the Spearman correlation  
 611 between these feature rankings. While the correlations degrade as we select more features (which is  
 612 to be expected), the marginal gains become similar among the remaining features after removing the  
 613 most important features.

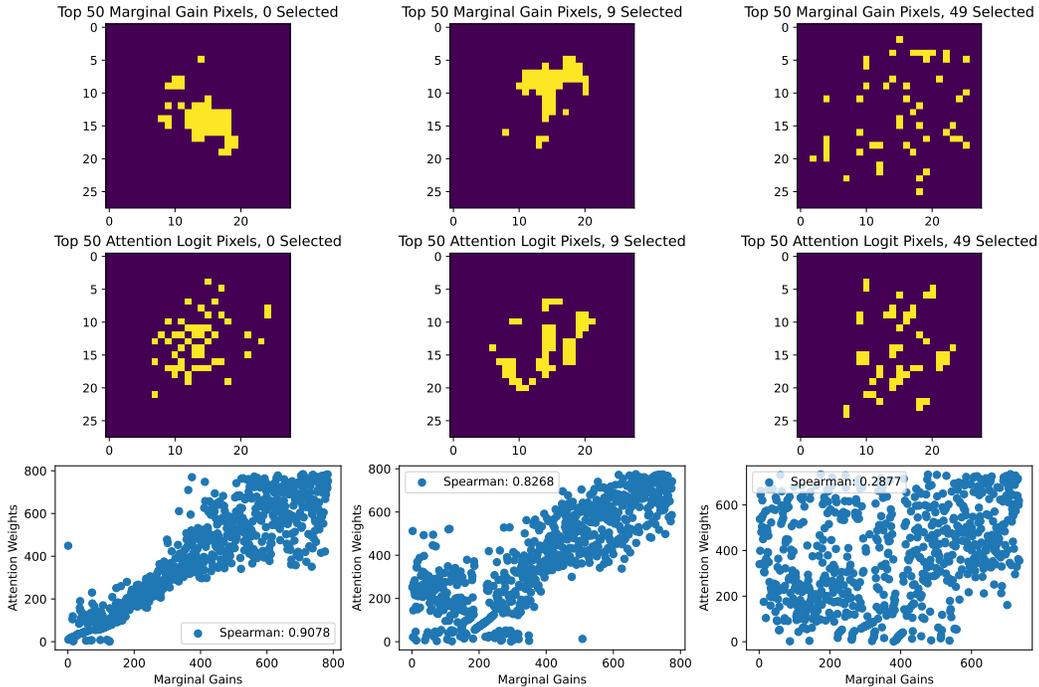


Figure 11: Marginal gain experiments. The first and second rows show that top 50 features chosen using the true marginal gains (top) and Sequential Attention (middle). The bottom row shows the Spearman correlation between these two computed sets of scores.