
Atlas: Universal Function Approximator For Memory Retention—Supplementary Appendices

Anonymous Author(s)

Affiliation

Address

email

1 The supplementary appendices contain additional experiments (Section A) and proofs (Section B) to
2 substantiate the claims in the main body of the paper.

3 A Additional experiments

4 The results for different experiments are given. Highlights from experiment A were discussed in the
5 main paper, but the training and validation loss curves are only presented in the appendix.

6 Experimentation was performed with Python and TensorFlow. Experiment A was performed on a
7 personal laptop with a 7th generation i7 Intel processor and took a few hours to finish thirty trials.
8 The loss function chosen for training and evaluation is the mean absolute error (MAE). The training
9 data set and test set in all experiments had 10000 data points, sampled uniformly at random. Gaussian
10 noise with standard deviation = 0.1 was added to all training and test data target values. The test set
11 was also used as a validation set to quantify the test error during training. All models were trained
12 with a learning rate of 0.01 with the Adam optimizer. All models and experiments used batch sizes
13 of 100 during training.

14 To test memory retention, two tasks, presented to an Atlas model one after the other, were constructed.
15 The details of each task are given below.

16 **Task 1** The training and test sets were sampled uniformly from the Task 1 target function over
17 the domain $[0., 1.]^2$, with Gaussian noise added to the target values. The initial Atlas model was
18 instantiated as a two-variable function that maps to a one-dimensional output, with $r = 0$ and $M = 0$
19 such that it is a minimally expressive model. The model was evaluated and trained for 30 epochs.
20 After training, the Atlas model was expanded using the built-in methods, such that r is increased by
21 one, and M is increased by two: $r' = r + 1$ and $M' = M + 2$. This training-expansion process was
22 repeated four times. The output of the model is presented at the end of each expansion iteration. The
23 target functions for Task 1 in each experiment is labeled $Y(\vec{x})$.

24 **Task 2** The test sets were sampled uniformly from the Task 2 target function over the domain
25 $[0., 1.]^2$, with Gaussian noise added to the target values. The training sets were sampled uniformly
26 over the domain $[0.45, 0.55]^2$, and target values of zero with added Gaussian noise. All models were
27 trained for 6 epochs. The Task 2 target function in each experiment is given by:

$$Y'(\vec{x}) = \begin{cases} 0 & 0.45 < x_i < 0.55 \ \forall i = 1, 2, \dots \\ Y(\vec{x}) & \text{otherwise.} \end{cases}$$

28 Task 2 effectively tests if a model changes only where new data was presented, with off-target effects
29 leading to larger test MAE.

30 A.1 Experiment A

31 A.1.1 Task 1

32 The Task 1 target function of Experiment A is given as follows, where the radius is measured from
 33 the centre of the domain $[0.5, 0.5]$:

$$r = \sqrt{(x_1 - \frac{1}{2})^2 + (x_2 - \frac{1}{2})^2}$$

$$\theta = \tan^{-1} \left((x - \frac{1}{2})^2, (y - \frac{1}{2})^2 \right)$$

$$Y_A = Y_A(x_1, x_2) = \sin(30r + \theta) + 2$$

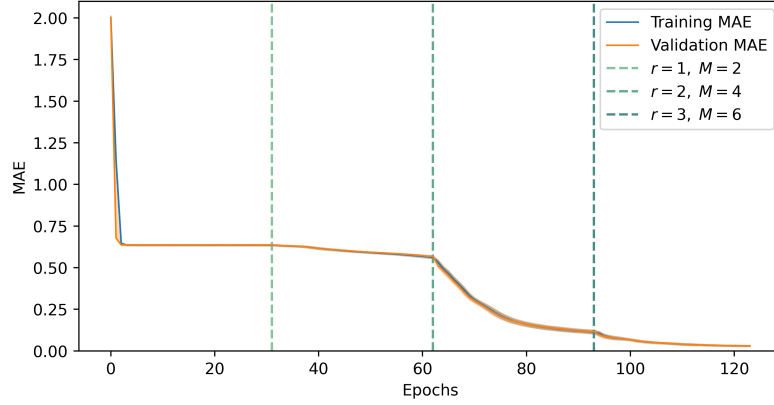


Figure 1: Training and validation MAE during the course of training on Task 1, Experiment A.

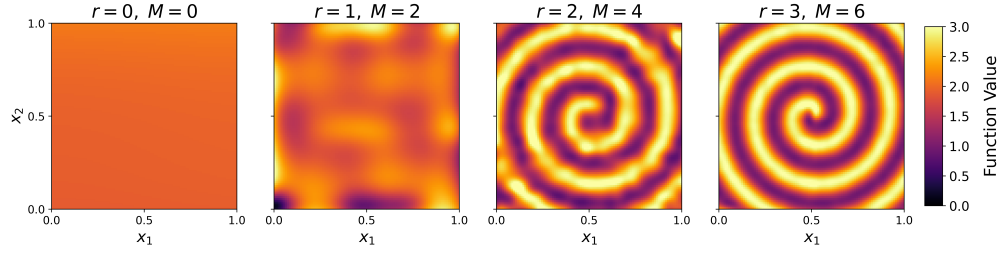


Figure 2: Outputs of the model during successive training and expansion iterations, Experiment A.

34 A.1.2 Task 2

35 The under-sampled target function Y'_A used for validation is given by:

$$Y'_A(x_1, x_2) = \begin{cases} 0 & 0.45 < x_i < 0.55 \ \forall i = 1, 2, \dots \\ Y_A(x_1, x_2) & \text{otherwise.} \end{cases}$$

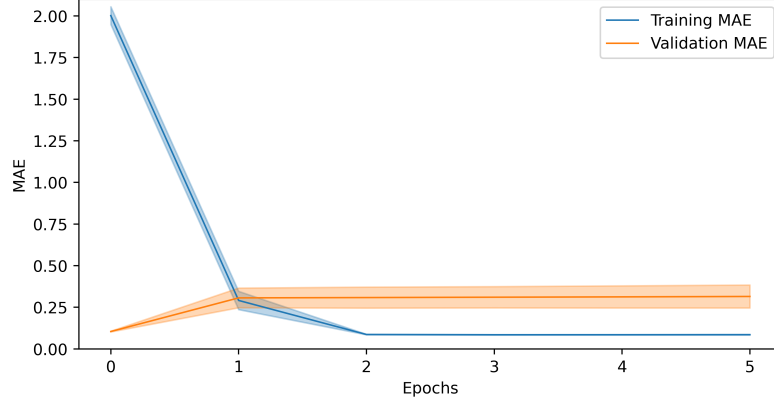


Figure 3: Training and validation MAE during the course of training on Task 2, Experiment A.

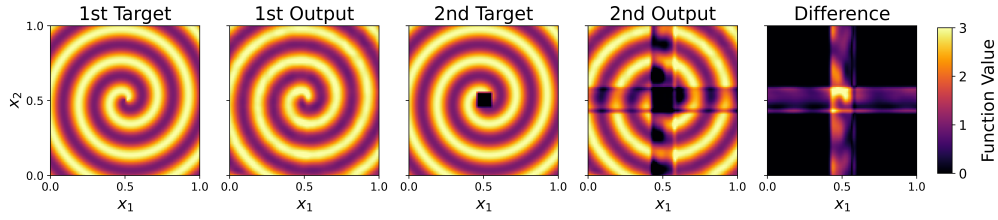


Figure 4: Visual inspection of target functions and model outputs over Task 1 and Task 2, Experiment A.

36 A.2 Experiment B

37 A.2.1 Task 1

38 The Task 1 target function for experiment B is given by:

$$Y_B(x_1, x_2) = \cos^2(20x_1 - 10) + \cos^2(10x_2 - 5) + \exp(-(20x_1 - 10)^2 - (20x_2 - 10)^2)$$

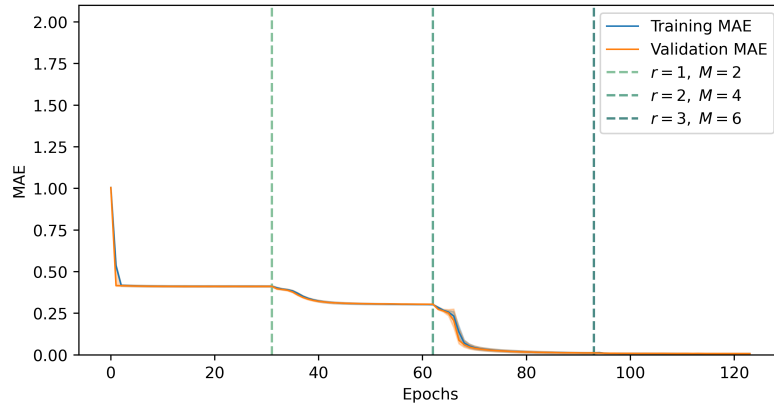


Figure 5: Training and validation MAE during the course of training on Task 1, Experiment B.

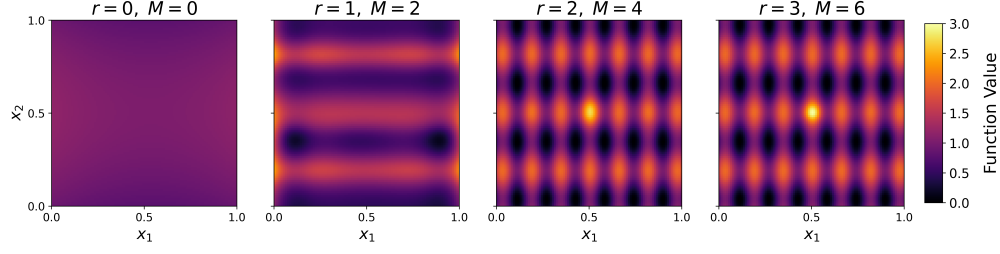


Figure 6: Outputs of the model during successive training and expansion iterations, Experiment B.

39 A.2.2 Task 2

40 The under-sampled target function Y'_B used for validation is given by:

$$Y'_B(x_1, x_2) = \begin{cases} 0 & 0.45 < x_i < 0.55 \forall i = 1, 2, \dots \\ Y_B(x_1, x_2) & \text{otherwise.} \end{cases}$$

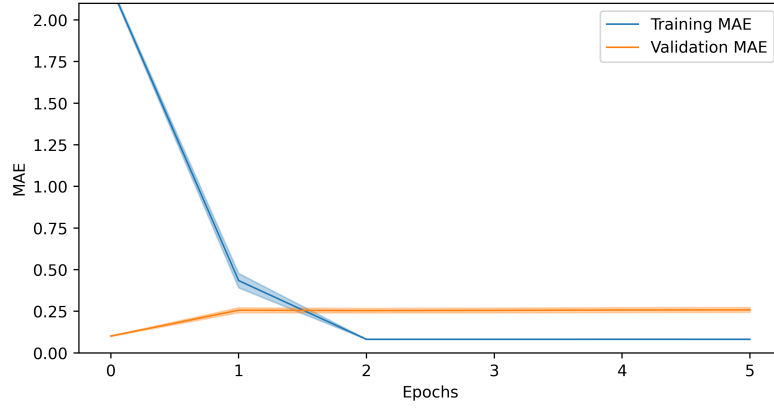


Figure 7: Training and validation MAE during the course of training on Task 2, Experiment B

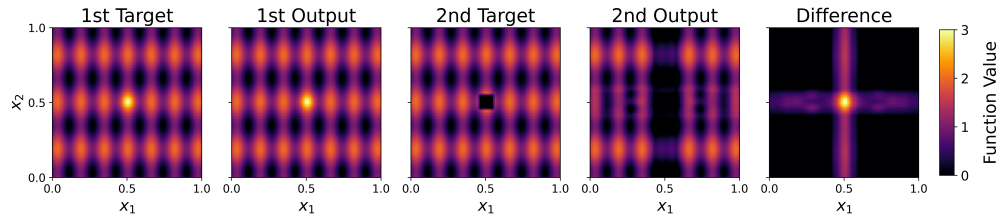


Figure 8: Visual inspection of target functions and model outputs over Task 1 and Task 2, Experiment B.

41 A.3 Experiment C

42 A.3.1 Task 1

43 The Task 1 target function for Experiment C is given by:

$$Y_C(x_1, x_2) = 2 + \cos(20x_1 - 10) \cos(20x_2 - 10)$$

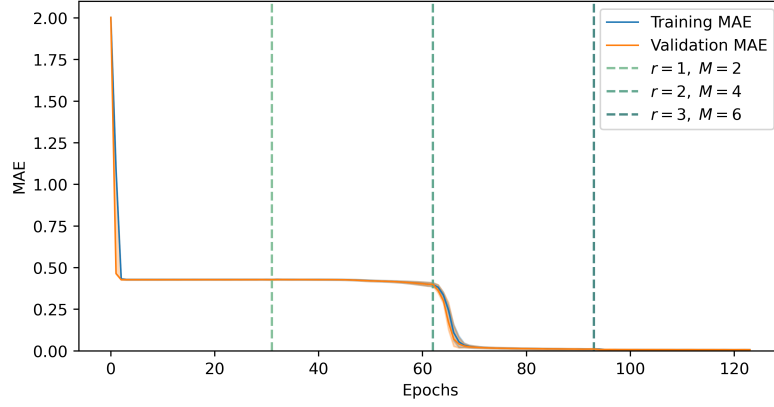


Figure 9: Training and validation MAE during the course of training on Task 1, Experiment C.

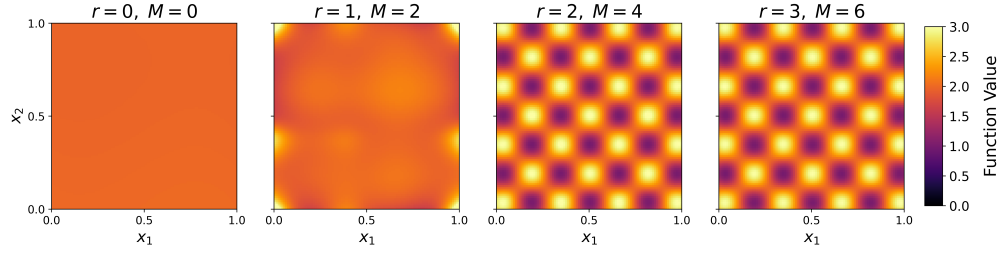


Figure 10: Outputs of the model during successive training and expansion iterations, Experiment C.

44 A.3.2 Task 2

45 The under-sampled target function Y'_C used for validation is given by:

$$Y'_C(x_1, x_2) = \begin{cases} 0 & 0.45 < x_i < 0.55 \forall i = 1, 2, \dots \\ Y_C(x_1, x_2) & \text{otherwise.} \end{cases}$$

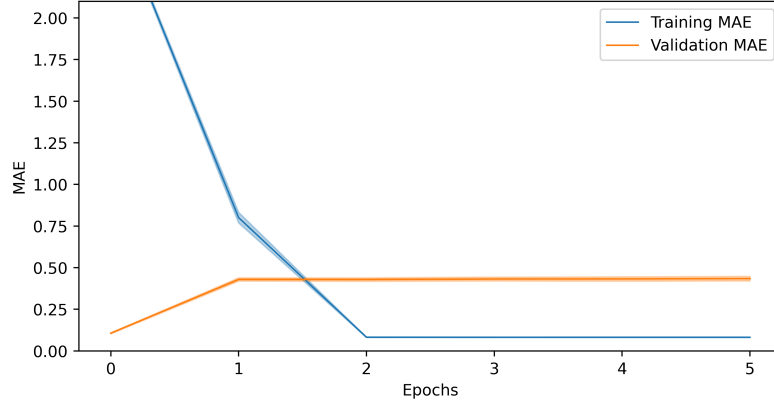


Figure 11: Training and validation MAE during the course of training on Task 2, Experiment C.

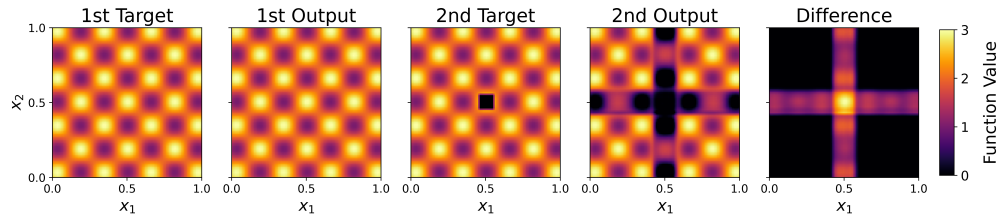


Figure 12: Visual inspection of target functions and model outputs over Task 1 and Task 2, Experiment C.

46 A.4 Experiment D

47 A.4.1 Task 1

48 The task 1 target function for experiment D is given by:

$$Y_D(x_1, x_2) = 2 + \sigma(\sin(2\pi x_1) \sin(2\pi x_2))$$

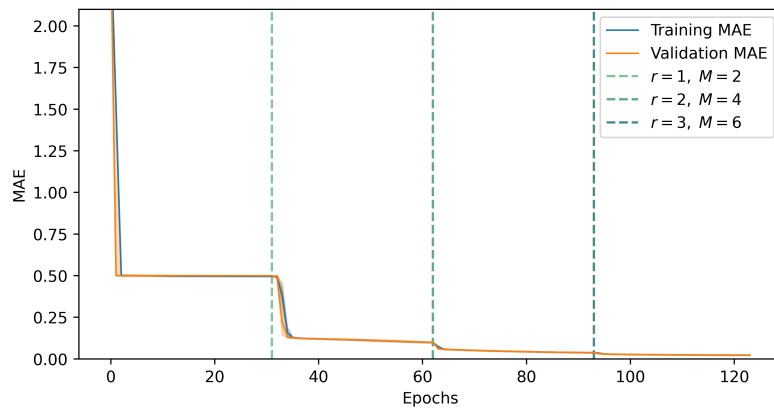


Figure 13: Training and validation MAE during the course of training on Task 1, Experiment D.

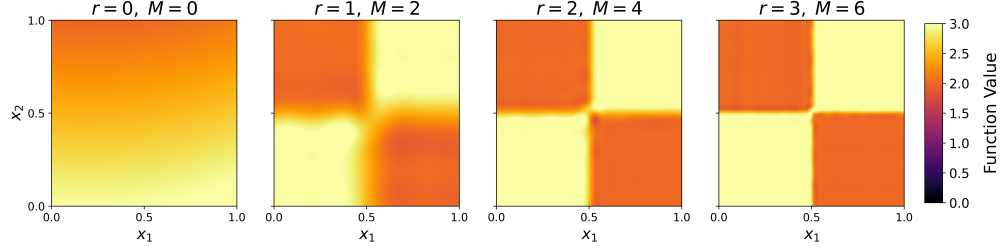


Figure 14: Outputs of the model during successive training and expansion iterations, Experiment D.

49 A.4.2 Task 2

50 The under-sampled target function Y'_D used for validation is given by:

$$Y'_D(x_1, x_2) = \begin{cases} 0 & 0.45 < x_i < 0.55 \forall i = 1, 2, \dots \\ Y_D(x_1, x_2) & \text{otherwise.} \end{cases}$$

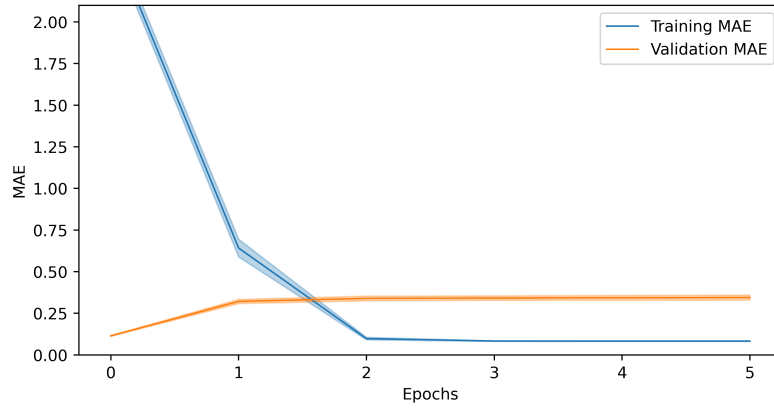


Figure 15: Training and validation MAE during the course of training on Task 2, Experiment D.

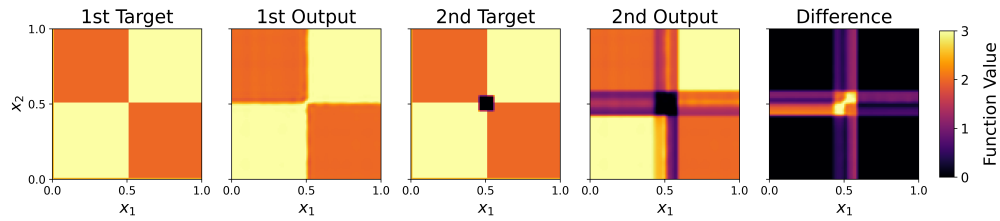


Figure 16: Visual inspection of target functions and model outputs over Task 1 and Task 2, Experiment D.

51 B Analysis and mathematical proofs

52 List of all definitions, theorems, corollaries, properties and their proofs are presented for completeness.
 53 The numbering of all statements match the main body of the paper. Some of the important results are
 54 also presented in the main body of the paper.

55 B.1 Stone-Weierstrass Theorem

56 Any continuous multi-variable function on a compact space can be uniformly approximated with
 57 multi-variable polynomials by the Stone-Weierstrass Theorem. Let \mathcal{I} denote an index set of tuples of

58 natural numbers including zero such that $i_j \in \mathbb{N}^0$ for all $j \in \mathbb{N}$ with $i = (i_1, \dots, i_n) \in \mathcal{I}$ and $a_i \in \mathbb{R}$.
 59 Multi-variable polynomials can be represented as:

$$y(\vec{x}) = y(x_1, \dots, x_n) = \sum_{i \in \mathcal{I}} a_i x_1^{i_1} x_2^{i_2} \dots x_n^{i_n} = \sum_{i \in \mathcal{I}} a_i \prod_{j=1}^n x_j^{i_j}$$

60 Each monomial term $a_i \prod_{j=1}^n x_j^{i_j}$ is a product of single-variable functions in each variable.

61 B.2 Exponential Representation Theorem

62 **Lemma 1.** *For any $a_i \in \mathbb{R}$, there exists $\gamma_i > 0$ and $\beta_i > 0$, such that: $a_i = \gamma_i - \beta_i$*

63 *Proof.* Let $a_i \in \mathbb{R}$. Three cases are considered.

64 If $a_i = 0$, then choose $\gamma_i = 1 > 0$ and $\beta_i = 1 > 0$, such that: $\gamma_i - \beta_i = 1 - 1 = 0 = a_i$

65 If $a_i > 0$, then choose $\gamma_i = a_i + 1 > 0$ and $\beta_i = 1 > 0$, such that: $\gamma_i - \beta_i = a_i + 1 - 1 = a_i$

66 If $a_i < 0$, then choose $\gamma_i = 1 > 0$ and $\beta_i = 1 + |a_i| > 0$, such that:

$$\gamma_i - \beta_i = 1 - (1 + |a_i|) = 1 - 1 - |a_i| = a_i$$

67

□

68 **Theorem 1** (Exponential representation theorem). *Any multi-variable polynomial function $p(\vec{x})$*
 69 *of n variables over the positive orthant, can be exactly represented by continuous single-variable*
 70 *functions $g_{i,j}(x_j)$ and $h_{i,j}(x_j)$ in the form:*

$$p(\vec{x}) = \sum_{i \in \mathcal{I}} \exp(\sum_{j=1}^n g_{i,j}(x_j)) - \exp(\sum_{j=1}^n h_{i,j}(x_j))$$

71 *Proof.* Consider any monomial term $a_i \prod_{j=1}^n x_j^{i_j}$ with $a_i \in \mathbb{R}$, then by Lemma 1 there exist strictly
 72 positive numbers $\gamma_i > 0$ and $\beta_i > 0$, such that:

$$\begin{aligned} a_i \prod_{j=1}^n x_j^{i_j} &= \gamma_i \prod_{j=1}^n x_j^{i_j} - \beta_i \prod_{j=1}^n x_j^{i_j} \\ &= \exp\left(\log\left(\gamma_i \prod_{j=1}^n x_j^{i_j}\right)\right) - \exp\left(\log\left(\beta_i \prod_{j=1}^n x_j^{i_j}\right)\right) \\ &= \exp\left(\log(\gamma_i) + \sum_{j=1}^n \log\left(x_j^{i_j}\right)\right) - \exp\left(\log(\beta_i) + \sum_{j=1}^n \log\left(x_j^{i_j}\right)\right) \end{aligned}$$

73 The argument of each exponential function is a sum of single-variable functions and constants.
 74 Without loss of generality, a set of single-variable functions can be defined such that:

$$a_i \prod_{j=1}^n x_j^{i_j} = \exp(\sum_{j=1}^n g_{i,j}(x_j)) - \exp(\sum_{j=1}^n h_{i,j}(x_j))$$

75 Since this holds for any $a_i \prod_{j=1}^n x_j^{i_j}$ and all $i \in \mathcal{I}$, it follows that:

$$p(\vec{x}) = \sum_{i \in \mathcal{I}} \exp(\sum_{j=1}^n g_{i,j}(x_j)) - \exp(\sum_{j=1}^n h_{i,j}(x_j))$$

76

□

77 There is a duality between representation and approximation. If any multi-variable polynomial can
 78 be exactly represented, then any continuous multi-variable function can be approximated to arbitrary
 79 accuracy.

80 B.2.1 Exponential approximation corollary

81 **Corollary 1** (Exponential approximation). *For any $\varepsilon > 0$, and continuous multi-variable function*
 82 *$y(\vec{x})$ of n variables over a compact domain in the positive orthant, there exist continuous single-*
 83 *variable functions $g_{i,j}(x_j)$ and $h_{i,j}(x_j)$ such that:*

$$\left| y(\vec{x}) - \left(\sum_{i \in \mathcal{I}} \exp(\sum_{j=1}^n g_{i,j}(x_j)) - \exp(\sum_{j=1}^n h_{i,j}(x_j)) \right) \right| < \varepsilon$$

84 *Proof.* Fix $\varepsilon > 0$, and let $y(\vec{x})$ be a continuous multi-variable function of n variables over a compact
 85 domain in the positive orthant.

86 By the Stone–Weierstrass theorem there exists a multi-variable polynomial $p(\vec{x})$ over the domain of
 87 $y(\vec{x})$ such that:

$$|y(\vec{x}) - p(\vec{x})| < \varepsilon$$

88 By Theorem 1, for any polynomial $p(\vec{x})$ over the positive orthant, there exist continuous single-
 89 variable functions $g_{i,j}(x_j)$ and $h_{i,j}(x_j)$ such that:

$$p(\vec{x}) = \sum_{i \in \mathcal{I}} \exp(\sum_{j=1}^n g_{i,j}(x_j)) - \exp(\sum_{j=1}^n h_{i,j}(x_j))$$

90 It follows that:

$$\left| y(\vec{x}) - \left(\sum_{i \in \mathcal{I}} \exp(\sum_{j=1}^n g_{i,j}(x_j)) - \exp(\sum_{j=1}^n h_{i,j}(x_j)) \right) \right| < \varepsilon$$

91

□

92 B.3 Single-variable function approximators

93 Each basis function S_i for a uniform cubic B-spline can be obtained by scaling and translating the
 94 input of the same activation function. The activation function is denoted $S(x)$ and is given by:

$$S(x) = \begin{cases} \frac{1}{6}x^3 & 0 \leq x < 1 \\ \frac{1}{6}[-3(x-1)^3 + 3(x-1)^2 + 3(x-1) + 1] & 1 \leq x < 2 \\ \frac{1}{6}[3(x-2)^3 - 6(x-2)^2 + 4] & 2 \leq x < 3 \\ \frac{1}{6}(4-x)^3 & 3 \leq x < 4 \\ 0 & \text{otherwise} \end{cases}$$

95 B.3.1 Definition of ρ -density B-spline functions

96 **Definition 1** (ρ -density B-spline function). A ρ -density B-spline function is a uniform cubic B-spline
 97 function with $2^{\rho+2}$ basis functions:

$$f(x) = \sum_{i=1}^{2^{\rho+2}} \theta_i S_i(x) = \sum_{i=1}^{2^{\rho+2}} \theta_i S(w_i x + b_i) = \sum_{i=1}^{2^{\rho+2}} \theta_i S((2^{\rho+2} - 3)x + 4 - i)$$

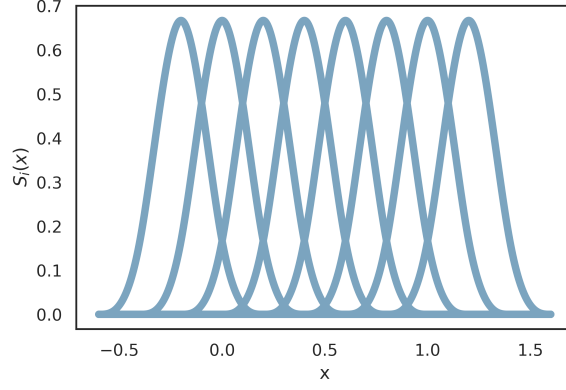


Figure 17: Set of eight uniform cubic B-spline basis functions, where $S_i(x) = S(w_i x + b_i)$.

98 B.3.2 Definition of mixed-density B-spline function

99 **Definition 2** (mixed-density B-spline function). A mixed-density B-spline function is a single-
 100 variable function approximator that is obtained by summing together different ρ -density B-spline
 101 functions. Only the maximum ρ -density B-spline function has trainable parameters, the others are
 102 constant. Mixed-density B-spline functions are of the form:

$$f(x) = \sum_{\rho=0}^r \sum_{i=1}^{2^{\rho+2}} \theta_{\rho,i} S_{\rho,i}(x)$$

103 The maximum density parameters $\theta_{r,i}$ are trainable, but the lower density parameters $\theta_{\rho,i}$ (with
 104 $\rho < r$) are in general non-zero constants. The function approximator can be expanded without losing
 105 previously learned values. Analytically, we can choose all the new parameters $\theta_{r+1,i} = 0, \forall i \in \mathbb{N}$
 106 such that:

$$f(x) = \sum_{\rho=0}^r \sum_{i=1}^{2^{\rho+2}} \theta_{\rho,i} S_{\rho,i}(x) = \sum_{\rho=0}^{r+1} \sum_{i=1}^{2^{\rho+2}} \theta_{\rho,i} S_{\rho,i}(x)$$

107 B.4 Atlas architecture

108 B.4.1 Atlas representation theorem

109 **Theorem 2** (Atlas representation theorem). Any multi-variable polynomial $p(\vec{x})$ of n variables
 110 over the positive orthant, can be exactly represented by continuous single-variable functions $f_j(x_j)$,
 111 $g_{i,j}(x_j)$, and $h_{i,j}(x_j)$ in the form:

$$p(\vec{x}) = \sum_{j=1}^n f_j(x_j) + \sum_{k=1}^{\infty} \frac{1}{k^2} \exp(\sum_{j=1}^n g_{k,j}(x_j)) - \frac{1}{k^2} \exp(\sum_{j=1}^n h_{k,j}(x_j))$$

112 *Proof.* Let $p(\vec{x})$ be a multi-variable polynomial over the positive orthant:

$$p(\vec{x}) = p(x_1, \dots, x_n) = \sum_{i \in \mathcal{I}} a_i x_1^{i_1} x_2^{i_2} \dots x_n^{i_n} = \sum_{i \in \mathcal{I}} a_i \prod_{j=1}^n x_j^{i_j}$$

113 Consider the set of terms that depend on at most one input variable, or single-variable terms in the
 114 expression for the polynomial $p(\vec{x})$:

$$\mathcal{P}_1 := \{a_i \prod_{j=1}^n x_j^{i_j} \mid i \in \mathcal{I}, i_j \neq 0 \implies i_k = 0, \forall k \neq j\}$$

115 It is worth noting that \mathcal{P}_1 contains the constant function.

116 Let \mathcal{Q} denote the index set of all single-variable monomial terms:

$$\mathcal{Q} := \{ i \mid i \in \mathcal{I}, i_j \neq 0 \implies i_k = 0, \forall k \neq j \}$$

117 The polynomial $p(\vec{x})$ can be rewritten in terms of single-variable functions and a residual polynomial
118 function p_{res} as:

$$\begin{aligned} p(\vec{x}) &= \sum_{i \in \mathcal{Q}} a_i \Pi_{j=1}^n x_j^{i_j} + \sum_{i \in \mathcal{I} \setminus \mathcal{Q}} a_i \Pi_{j=1}^n x_j^{i_j} \\ &= \sum_{j=1}^n f_j(x_j) + p_{res}(\vec{x}) \end{aligned}$$

119 The single-variable terms can be consumed by a sum of n arbitrary single-variable functions $f_j(x_j)$.

120 By Theorem 1, for any polynomial $p_{res}(\vec{x})$ over the positive orthant, there exist continuous single-
121 variable functions $g_{i,j}(x_j)$ and $h_{i,j}(x_j)$ such that:

$$p_{res}(\vec{x}) = \sum_{i \in \mathcal{I} \setminus \mathcal{Q}} \exp(\sum_{j=1}^n g_{i,j}(x_j)) - \exp(\sum_{j=1}^n h_{i,j}(x_j))$$

122 Since the index set is countable, one can use another indexing scheme:

$$p_{res}(\vec{x}) = \sum_{k=1}^{\infty} \exp(\sum_{j=1}^n g_{k,j}(x_j)) - \exp(\sum_{j=1}^n h_{k,j}(x_j))$$

123 Scale factors can be introduced without changing the representation:

$$\begin{aligned} p_{res}(\vec{x}) &= \sum_{k=1}^{\infty} \exp(\log k^2 - \log k^2 + \sum_{j=1}^n g_{k,j}(x_j)) - \exp(\log k^2 - \log k^2 + \sum_{j=1}^n h_{k,j}(x_j)) \\ &= \sum_{k=1}^{\infty} \frac{1}{k^2} \exp(\log k^2 + \sum_{j=1}^n g_{k,j}(x_j)) - \frac{1}{k^2} \exp(\log k^2 + \sum_{j=1}^n h_{k,j}(x_j)) \end{aligned}$$

124 Since the single-variable functions $g_{i,j}(x_j)$ and $h_{i,j}(x_j)$ are arbitrary, one can absorb the constants
125 and redefine $g_{i,j}(x_j)$ and $h_{i,j}(x_j)$ to obtain:

$$p_{res}(\vec{x}) = \sum_{k=1}^{\infty} \frac{1}{k^2} \exp(\sum_{j=1}^n g_{k,j}(x_j)) - \frac{1}{k^2} \exp(\sum_{j=1}^n h_{k,j}(x_j))$$

126 Substituting the expressions one obtains:

$$p(\vec{x}) = \sum_{j=1}^n f_j(x_j) + \sum_{k=1}^{\infty} \frac{1}{k^2} \exp(\sum_{j=1}^n g_{k,j}(x_j)) - \frac{1}{k^2} \exp(\sum_{j=1}^n h_{k,j}(x_j))$$

127

□

128 There is a duality between representation and approximation. If any multi-variable polynomial can
129 be exactly represented, then any continuous multi-variable function can be approximated to arbitrary
130 accuracy.

131 B.4.2 Atlas definition

132 **Definition 3** (Atlas). Atlas is a function approximator of n variables, with mixed-density B-spline
 133 functions $f_j(x_j)$, $g_{i,j}(x_j)$, and $h_{i,j}(x_j)$ in the form:

$$A(\vec{x}) := \sum_{j=1}^n f_j(x_j) + \sum_{k=1}^M \frac{1}{k^2} \exp(\sum_{j=1}^n g_{k,j}(x_j)) - \frac{1}{k^2} \exp(\sum_{j=1}^n h_{k,j}(x_j))$$

134 Atlas is equivalently given by the compact notation:

$$\begin{aligned} A(\vec{x}) &:= \sum_{j=1}^n f_j(x_j) + \sum_{k=1}^M \frac{1}{k^2} \exp(\sum_{j=1}^n g_{k,j}(x_j)) - \frac{1}{k^2} \exp(\sum_{j=1}^n h_{k,j}(x_j)) \\ &= F(\vec{x}) + \sum_{k=1}^M \frac{1}{k^2} \exp(G_k(\vec{x})) - \frac{1}{k^2} \exp(H_k(\vec{x})) \\ &= F(\vec{x}) + G(\vec{x}) - H(\vec{x}) \end{aligned}$$

135 B.4.3 Atlas polynomial approximation

136 **Theorem 3** (Atlas polynomial approximation). *For any multi-variable polynomial $p(\vec{x})$ over the*
 137 *positive orthant with bounded and compact domain $D(p)$ and $\varepsilon > 0$, there exists an Atlas model*
 138 *$A(\vec{x})$ such that:*

$$|p(\vec{x}) - A(\vec{x})| < \varepsilon$$

139 *Proof.* Let $p(\vec{x})$ be a multi-variable polynomial $p(\vec{x})$ of n variables over the positive orthant, and fix
 140 $\varepsilon > 0$, and choose:

$$\varepsilon = \varepsilon_1 + \varepsilon_2$$

141 By Theorem 2, there exist continuous single-variable functions $f_j(x_j)$, $g_{i,j}(x_j)$, and $h_{i,j}(x_j)$ such
 142 that:

$$p(\vec{x}) = \sum_{j=1}^n f_j(x_j) + \sum_{k=1}^{\infty} \frac{1}{k^2} \exp(\sum_{j=1}^n g_{k,j}(x_j)) - \frac{1}{k^2} \exp(\sum_{j=1}^n h_{k,j}(x_j))$$

143 If $p(\vec{x})$ has finitely many terms, then let M denote the number of residual terms:

$$\begin{aligned} p(\vec{x}) &= \sum_{j=1}^n f_j(x_j) + \sum_{k=1}^M \frac{1}{k^2} \exp(\sum_{j=1}^n g_{k,j}(x_j)) - \frac{1}{k^2} \exp(\sum_{j=1}^n h_{k,j}(x_j)) \\ &= F(\vec{x}) + \sum_{k=1}^M \frac{1}{k^2} \exp(G_k(\vec{x})) - \frac{1}{k^2} \exp(H_k(\vec{x})) \\ &= F(\vec{x}) + G(\vec{x}) - H(\vec{x}) \end{aligned}$$

144 Choose mixed-density B-spline functions $f_j^*(x_j)$, $g_{i,j}^*(x_j)$, and $h_{i,j}^*(x_j)$ such that the Atlas model
 145 $A(\vec{x})$ is given by:

$$\begin{aligned}
A^*(\vec{x}) &= \sum_{j=1}^n f_j^*(x_j) + \sum_{k=1}^M \frac{1}{k^2} \exp(\sum_{j=1}^n g_{k,j}^*(x_j)) - \frac{1}{k^2} \exp(\sum_{j=1}^n h_{k,j}^*(x_j)) \\
&= F^*(\vec{x}) + \sum_{k=1}^M \frac{1}{k^2} \exp(G_k^*(\vec{x})) - \frac{1}{k^2} \exp(H_k^*(\vec{x})) \\
&= F^*(\vec{x}) + G^*(\vec{x}) - H^*(\vec{x})
\end{aligned}$$

146 Then it follows that,

$$\begin{aligned}
|p(\vec{x}) - A(\vec{x})| &= |F(\vec{x}) + G(\vec{x}) - H(\vec{x}) - (F^*(\vec{x}) + G^*(\vec{x}) - H^*(\vec{x}))| \\
&= |F(\vec{x}) - F^*(\vec{x}) + G(\vec{x}) - G^*(\vec{x}) - (H(\vec{x}) - H^*(\vec{x}))| \\
&\leq |F(\vec{x}) - F^*(\vec{x})| + |G(\vec{x}) - G^*(\vec{x})| + |H(\vec{x}) - H^*(\vec{x})|
\end{aligned}$$

147 The first set of functions is easily shown to have bounded error. Choose mixed-density B-spline
148 functions $f_j^*(x_j)$ such that:

$$|f_j(x_j) - f_j^*(x_j)| < \frac{\varepsilon_1}{n}$$

149 Then it follows,

$$\begin{aligned}
|F(\vec{x}) - F^*(\vec{x})| &= \left| \sum_{j=1}^n f_j(x_j) - \sum_{j=1}^n f_j^*(x_j) \right| \\
&\leq \sum_{j=1}^n |f_j(x_j) - f_j^*(x_j)| \\
&< \varepsilon_1
\end{aligned}$$

150 The interior functions for the exponential functions are more complicated.

151 *Remark.* The uniform continuity of exponentials on bounded domains makes it possible to bound the
152 approximation error in each exponential term. The exponential function is uniformly continuous
153 on a compact and bounded subset of the real numbers $[a, b]$. Thus, for any $\varepsilon_{\text{exp}} > 0$, there exists a
154 $\delta_{\text{exp}} > 0$, such that for every $x, y \in [a, b]$:

$$|x - y| < \delta_{\text{exp}} \implies |\exp(x) - \exp(y)| < \varepsilon_{\text{exp}}$$

155 For all exponential functions on bounded and compact domains choose:

$$\varepsilon_{\text{exp}} = \frac{3\varepsilon_2}{\pi^2}$$

156 Choose the smallest δ_{exp} for all M exponential functions, so that the implication holds. Choose
157 $\delta_{g,k,j}$, such that:

$$\sum_{j=1}^n \delta_{g,k,j} < \delta_{\text{exp}}$$

158 Choose mixed-density B-spline functions $g_{i,j}^*(x_j)$, and $h_{i,j}^*(x_j)$ such that:

$$\begin{aligned} |g_{k,j}(x_j) - g_{k,j}^*(x_j)| &< \delta_{g,k,j} \\ |h_{k,j}(x_j) - h_{k,j}^*(x_j)| &< \delta_{g,k,j} \end{aligned}$$

159 The interior functions $g_{k,j}^*(x_j)$ have bounded approximation error $\delta_{g,k,j}$ one obtains:

$$\begin{aligned} &\left| \sum_{j=1}^n g_{k,j}(x_j) - \sum_{j=1}^n g_{k,j}^*(x_j) \right| \\ &\leq \sum_{j=1}^n |g_{k,j}(x_j) - g_{k,j}^*(x_j)| \\ &< \sum_{j=1}^n \delta_{g,k,j} < \delta_{\text{exp}} \end{aligned}$$

160 This implies that:

$$\left| \exp\left(\sum_{j=1}^n g_{k,j}(x_j)\right) - \exp\left(\sum_{j=1}^n g_{k,j}^*(x_j)\right) \right| < \varepsilon_{\text{exp}}$$

161 Recombining this result with the exponential terms yields:

$$\begin{aligned} |G(\vec{x}) - G^*(\vec{x})| &= \left| \sum_{k=1}^M \frac{1}{k^2} \exp\left(\sum_{j=1}^n g_{k,j}(x_j)\right) - \sum_{k=1}^M \frac{1}{k^2} \exp\left(\sum_{j=1}^n g_{k,j}^*(x_j)\right) \right| \\ |G(\vec{x}) - G^*(\vec{x})| &\leq \sum_{k=1}^M \frac{1}{k^2} \left| \exp\left(\sum_{j=1}^n g_{k,j}(x_j)\right) - \exp\left(\sum_{j=1}^n g_{k,j}^*(x_j)\right) \right| \\ |G(\vec{x}) - G^*(\vec{x})| &< \sum_{k=1}^M \frac{1}{k^2} \varepsilon_{\text{exp}} \end{aligned}$$

162 The scaling factors of k^{-2} were chosen for convergence, such that:

$$\begin{aligned} |G(\vec{x}) - G^*(\vec{x})| &< \sum_{k=1}^{\infty} \frac{1}{k^2} \varepsilon_{\text{exp}} \\ |G(\vec{x}) - G^*(\vec{x})| &< \varepsilon_{\text{exp}} \sum_{k=1}^{\infty} \frac{1}{k^2} \\ |G(\vec{x}) - G^*(\vec{x})| &< \varepsilon_{\text{exp}} \frac{\pi^2}{6} = \frac{3\varepsilon_2}{\pi^2} \frac{\pi^2}{6} = \frac{\varepsilon_2}{2} \end{aligned}$$

163 It follows that:

$$|G(\vec{x}) - G^*(\vec{x})| < \frac{\varepsilon_2}{2}$$

164 The same argument holds for $|H(\vec{x}) - H^*(\vec{x})|$, and one obtains the result:

$$\begin{aligned} |p(\vec{x}) - A(\vec{x})| &\leq |F(\vec{x}) - F^*(\vec{x})| + |G(\vec{x}) - G^*(\vec{x})| + |H(\vec{x}) - H^*(\vec{x})| \\ |p(\vec{x}) - A(\vec{x})| &< \varepsilon_1 + \frac{\varepsilon_2}{2} + \frac{\varepsilon_2}{2} \\ |p(\vec{x}) - A(\vec{x})| &< \varepsilon_1 + \varepsilon_2 \end{aligned}$$

165 Finally,

$$|p(\vec{x}) - A(\vec{x})| < \varepsilon$$

166

□

167 **B.4.4 Universal function approximation theorem**

168 **Theorem 4** (Atlas universal function approximation). *For any $\varepsilon > 0$, and continuous multi-variable*
 169 *function $y(\vec{x})$ of n variables over a compact domain in the positive orthant, there exists an Atlas*
 170 *model $A(\vec{x})$ such that:*

$$|y(\vec{x}) - A(\vec{x})| < \varepsilon$$

171 *Proof.* Let $\varepsilon > 0$, and $y(\vec{x})$ be a continuous multi-variable function of n variables over a compact
 172 domain in the positive orthant. Choose $\varepsilon_1 + \varepsilon_2 = \varepsilon$.

173 By the Stone–Weierstrass theorem there exists a multi-variable polynomial $p(\vec{x})$ over the domain of
 174 $y(\vec{x})$ such that:

$$|y(\vec{x}) - p(\vec{x})| < \varepsilon_1$$

175 By Theorem 3 an Atlas model can approximate the polynomial $p(\vec{x})$ to arbitrary precision:

$$|p(\vec{x}) - A(\vec{x})| < \varepsilon_2$$

176 It follows from the triangle inequality that:

$$\begin{aligned} |y(\vec{x}) - A(\vec{x})| &\leq |y(\vec{x}) - p(\vec{x})| + |p(\vec{x}) - A(\vec{x})| \\ |y(\vec{x}) - A(\vec{x})| &< \varepsilon_1 + \varepsilon_2 \end{aligned}$$

177 Finally,

$$|y(\vec{x}) - A(\vec{x})| < \varepsilon$$

178

□

179 **B.5 Atlas properties**

180 **B.5.1 Atlas expansion**

181 The number of exponential terms can be increased without changing the output of the model. We can
 182 choose to initialise $G_{M+1}(\vec{x}) = 0$ and $H_{M+1}(\vec{x}) = 0$, such that the model capacity can be increased
 183 without changing the output of the model:

$$\begin{aligned} A_{M+1}(\vec{x}) &= \sum_{k=1}^{M+1} \frac{1}{k^2} \exp(G_k(\vec{x})) - \frac{1}{k^2} \exp(H_k(\vec{x})) \\ &= \frac{1}{(M+1)^2} \exp(G_{M+1}(\vec{x})) - \frac{1}{(M+1)^2} \exp(H_{M+1}(\vec{x})) + A(\vec{x}) \\ &= \frac{1}{(M+1)^2} \exp(0) - \frac{1}{(M+1)^2} \exp(0) + A(\vec{x}) \\ &= A(\vec{x}) \end{aligned}$$

184 The density of basis functions in Atlas can also be incremented without changing the learned output
 185 of the model. The density of basis functions can also be increased without changing the output for
 186 any of mixed-density B-spline functions Ψ of the form:

$$\Psi(x) = \sum_{\rho=0}^r \sum_{i=1}^{2^{\rho+2}} \theta_{\rho,i} S_{\rho,i}(x)$$

187 Analytically, we can choose all the new parameters $\theta_{r+1,i} = 0, \forall i \in \mathbb{N}$ such that:

$$\Psi(x) = \sum_{\rho=0}^r \sum_{i=1}^{2^{\rho+2}} \theta_{\rho,i} S_{\rho,i}(x) = \sum_{\rho=0}^{r+1} \sum_{i=1}^{2^{\rho+2}} \theta_{\rho,i} S_{\rho,i}(x)$$

188 The last thing to note is that only the parameters for the largest specified density are trainable, in
 189 contrast to smaller density parameters that are fixed constants.

190 B.5.2 Atlas sparsity

191 **Property 1** (Sparsity). For any $\vec{x} \in D(A) \subset R^n$ and bounded trainable parameters θ_i with index set
 192 Θ , the gradient vector of trainable parameters for Atlas is sparse:

$$\left\| \vec{\nabla}_{\vec{\theta}} A(\vec{x}) \right\|_0 = \sum_{i \in \Theta} d_{\text{Hamming}} \left(\frac{\partial A}{\partial \theta_i}(\vec{x}), 0 \right) \leq 4n(2M+1)$$

193 *Proof.* Let $A(\vec{x})$ denote some Atlas model, with mixed-density B-spline functions $f_j(x_j), g_{i,j}(x_j),$
 194 and $h_{i,j}(x_j)$ in the form:

$$A(\vec{x}) = \sum_{j=1}^n f_j(x_j) + \sum_{k=1}^M \frac{1}{k^2} \exp(\sum_{j=1}^n g_{k,j}(x_j)) - \frac{1}{k^2} \exp(\sum_{j=1}^n h_{k,j}(x_j))$$

195 Each mixed-density B-spline function has its own parameters that are independent of every other
 196 mixed-density B-spline. The mixed-density B-splines function $\Psi(x)$ is by definition given by:

$$\Psi(x) = \sum_{\rho=0}^r \sum_{i=1}^{2^{\rho+2}} \theta_{\rho,i} S_{\rho,i}(x)$$

197 Thus for every mixed-density B-spline function in $A(\vec{x})$:

$$\begin{aligned} f_j(x_j) &= \sum_{\rho=0}^r \sum_{i=1}^{2^{\rho+2}} \theta_{f,(\rho,i,j)} S_{\rho,i}(x_j) \\ g_{k,j}(x_j) &= \sum_{\rho=0}^r \sum_{i=1}^{2^{\rho+2}} \theta_{g,(\rho,i,k,j)} S_{\rho,i}(x) \\ h_{k,j}(x_j) &= \sum_{\rho=0}^r \sum_{i=1}^{2^{\rho+2}} \theta_{h,(\rho,i,k,j)} S_{\rho,i}(x) \end{aligned}$$

198 Only the maximum density basis function $\rho = r$ have trainable parameters. The maximum density
 199 r -density B-spline functions are uniform B-spline functions with trainable parameters. There are at
 200 most four basis functions that are non-zero for any given x_j , and as such the gradient vector with
 201 respect to trainable parameters will have at most four non-zero entries for each r -density B-spline

202 function, the same four parameters for each mixed-density B-spline functions $f_j(x_j)$, $g_{i,j}(x_j)$, and
 203 $h_{i,j}(x_j)$. One simply needs to count the number of mixed-density B-spline functions.
 204 The number of mixed-density B-spline functions labeled $f_j(x_j)$ are in total n , with 4 active trainable
 205 parameters each.
 206 The number of functions labeled $g_{k,j}(x_j)$ are in total nM . For each M , there are n mixed-density
 207 B-spline functions, with 4 active trainable parameters each.
 208 The number of mixed-density B-spline functions labeled $h_{k,j}(x_j)$ are in total nM . For each M , there
 209 are n mixed-density B-spline functions, with 4 active trainable parameters each.
 210 The total number of active trainable parameters is thus:

$$4n + 4nM + 4nM = 4n(2M + 1)$$

211

□

212 *Remark.* The total number of trainable parameters for each mixed-density B-spline function is 2^{r+2} .
 213 For a fixed number of variables n , the model has a total of $2^{r+2}n(2M + 1)$ trainable parameters. The
 214 gradient vector has a maximum of $4n(2M + 1)$ non-zero entries, which is independent of r . Recall
 215 that only the maximum density ($\rho = r$) cubic B-spline function has trainable parameters. The fraction
 216 of trainable basis functions that are active is at most 2^{-r} . Sparsity entails efficient implementation,
 217 and suggests possible memory retention and robustness to catastrophic forgetting.

218 It is worth noting that the total number of parameters (including constants) is:

$$\text{Total number of parameters} \propto \sum_{\rho=0}^r 2^{\rho+2}n(2M + 1) \approx 2^{r+1}n(2M + 1)$$

219 B.5.3 Atlas gradient flow attenuation

220 **Property 2** (Gradient flow attenuation). For any $\vec{x} \in D(A) \subset R^n$ and bounded trainable parameters
 221 θ_i with index set Θ : if all the mixed-density B-spline functions are bounded, then the gradient vector
 222 of trainable parameters for Atlas is bounded:

$$\left\| \vec{\nabla}_{\vec{\theta}} A(\vec{x}) \right\|_1 = \sum_{i \in \Theta} \left| \frac{\partial A}{\partial \theta_i}(\vec{x}) \right| < U$$

223 *Proof.* Let $A(\vec{x})$ denote some Atlas model, with mixed-density B-spline functions $f_j(x_j)$, $g_{i,j}(x_j)$,
 224 and $h_{i,j}(x_j)$ in the form:

$$\begin{aligned} A(\vec{x}) &= \sum_{j=1}^n f_j(x_j) + \sum_{k=1}^M \frac{1}{k^2} \exp(\sum_{j=1}^n g_{k,j}(x_j)) - \frac{1}{k^2} \exp(\sum_{j=1}^n h_{k,j}(x_j)) \\ &= F(\vec{x}) + \sum_{k=1}^M \frac{1}{k^2} \exp(G_k(\vec{x})) - \frac{1}{k^2} \exp(H_k(\vec{x})) \\ &= F(\vec{x}) + G(\vec{x}) - H(\vec{x}) \end{aligned}$$

225 With each mixed-density B-spline function in $A(\vec{x})$ given by:

$$\begin{aligned}
f_j(x_j) &= \sum_{\rho=0}^r \sum_{i=1}^{2^{\rho+2}} \theta_{f,(\rho,i,j)} S_{\rho,i}(x_j) \\
g_{k,j}(x_j) &= \sum_{\rho=0}^r \sum_{i=1}^{2^{\rho+2}} \theta_{g,(\rho,i,k,j)} S_{\rho,i}(x) \\
h_{k,j}(x_j) &= \sum_{\rho=0}^r \sum_{i=1}^{2^{\rho+2}} \theta_{h,(\rho,i,k,j)} S_{\rho,i}(x)
\end{aligned}$$

226 The norm of the gradient of $A(\vec{x})$ with respect to trainable parameters is given by:

$$\begin{aligned}
\|\vec{\nabla}_{\vec{\theta}} A(\vec{x})\|_1 &= \|\vec{\nabla}_{\vec{\theta}} (F(\vec{x}) + G(\vec{x}) - H(\vec{x}))\|_1 \\
\|\vec{\nabla}_{\vec{\theta}} A(\vec{x})\|_1 &= \|\vec{\nabla}_{\vec{\theta}} F(\vec{x}) + \vec{\nabla}_{\vec{\theta}} G(\vec{x}) - \vec{\nabla}_{\vec{\theta}} H(\vec{x})\|_1 \\
\|\vec{\nabla}_{\vec{\theta}} A(\vec{x})\|_1 &\leq \|\vec{\nabla}_{\vec{\theta}} F(\vec{x})\|_1 + \|\vec{\nabla}_{\vec{\theta}} G(\vec{x})\|_1 + \|\vec{\nabla}_{\vec{\theta}} H(\vec{x})\|_1
\end{aligned}$$

227 The first term is bounded,

$$\begin{aligned}
\|\vec{\nabla}_{\vec{\theta}} F(\vec{x})\|_1 &= \left\| \vec{\nabla}_{\vec{\theta}} \left(\sum_{j=1}^n f_j(x_j) \right) \right\|_1 \\
&= \left\| \sum_{j=1}^n \vec{\nabla}_{\vec{\theta}} f_j(x_j) \right\|_1 \\
&\leq \sum_{j=1}^n \|\vec{\nabla}_{\vec{\theta}} f_j(x_j)\|_1
\end{aligned}$$

228 Substituting the expression for each $f_j(x_j)$, all lower densities have constant parameters:

$$\begin{aligned}
\|\vec{\nabla}_{\vec{\theta}} F(\vec{x})\|_1 &\leq \sum_{j=1}^n \left\| \vec{\nabla}_{\vec{\theta}} \left(\sum_{\rho=0}^r \sum_{i=1}^{2^{\rho+2}} \theta_{f,(\rho,i,j)} S_{\rho,i}(x_j) \right) \right\|_1 \\
\|\vec{\nabla}_{\vec{\theta}} F(\vec{x})\|_1 &\leq \sum_{j=1}^n \left\| \vec{\nabla}_{\vec{\theta}} \left(\sum_{i=1}^{2^{r+2}} \theta_{f,(r,i,j)} S_{r,i}(x_j) \right) \right\|_1 \\
\|\vec{\nabla}_{\vec{\theta}} F(\vec{x})\|_1 &\leq \sum_{j=1}^n \sum_{i=1}^{2^{r+2}} \|\vec{\nabla}_{\vec{\theta}} (\theta_{f,(r,i,j)} S_{r,i}(x_j))\|_1 \\
\|\vec{\nabla}_{\vec{\theta}} F(\vec{x})\|_1 &\leq \sum_{j=1}^n \sum_{i=1}^{2^{r+2}} |S_{r,i}(x_j)|
\end{aligned}$$

229 Each basis function is continuous and bounded by some positive constant $u > 0$, such that $S(x) < u$
230 regardless of its density, and it follows that:

$$\|\vec{\nabla}_{\vec{\theta}} F(\vec{x})\|_1 \leq \sum_{j=1}^n \sum_{i=1}^{2^{r+2}} u$$

231 The last thing to include is that at most four basis functions are non-zero, regardless of the value of r ,
 232 so a tighter upper bound is:

$$\left\| \vec{\nabla}_{\vec{\theta}} F(\vec{x}) \right\|_1 \leq \sum_{j=1}^n \sum_{i=1}^4 u = 4nu$$

233 *Remark.* Each ρ -density B-spline function has at most four active basis functions, and each mixed-
 234 density B-spline function has $r + 1$ different ρ -density B-spline functions. If the lower densities
 235 $\rho < r$ were also trainable, then this upper bound would instead be $4nu(r + 1)$. This is why only the
 236 maximum density was chosen to be trainable.

237 The exponential terms are more complicated.

$$\begin{aligned} \left\| \vec{\nabla}_{\vec{\theta}} G(\vec{x}) \right\|_1 &= \left\| \vec{\nabla}_{\vec{\theta}} \left(\sum_{k=1}^M \frac{1}{k^2} \exp(G_k(\vec{x})) \right) \right\|_1 \\ \left\| \vec{\nabla}_{\vec{\theta}} G(\vec{x}) \right\|_1 &\leq \sum_{k=1}^M \frac{1}{k^2} \left\| \vec{\nabla}_{\vec{\theta}} (\exp(G_k(\vec{x}))) \right\|_1 \\ \left\| \vec{\nabla}_{\vec{\theta}} G(\vec{x}) \right\|_1 &\leq \sum_{k=1}^M \frac{1}{k^2} \left\| \exp(G_k(\vec{x})) \vec{\nabla}_{\vec{\theta}} (G_k(\vec{x})) \right\|_1 \\ \left\| \vec{\nabla}_{\vec{\theta}} G(\vec{x}) \right\|_1 &\leq \sum_{k=1}^M \frac{1}{k^2} \exp(G_k(\vec{x})) \left\| \vec{\nabla}_{\vec{\theta}} (G_k(\vec{x})) \right\|_1 \end{aligned}$$

238 Each mixed-density B-spline function is bounded, so

$$|g_{k,j}(x_j)| = \left| \sum_{\rho=0}^r \sum_{i=1}^{2^{\rho+2}} \theta_{g,(\rho,i,k,j)} S_{\rho,i}(x) \right| < u_{g,(k,j)}$$

239 Since n is fixed and finite, the functions $G_k(\vec{x})$ are bounded:

$$|G_k(\vec{x})| = \left| \sum_{j=1}^n g_{k,j}(x_j) \right| \leq \sum_{j=1}^n |g_{k,j}(x_j)| < \sum_{j=1}^n u_{g,(k,j)} = u_{g,(k)}$$

240 Since this is true for each G_k , one can choose the maximum bound:

$$u_g = \max_{k=1,\dots,M} \{u_{g,(k)}\}$$

241 It is evident that:

$$G_k(\vec{x}) \leq |G_k(\vec{x})| < u_g$$

242 Since the exponential function is monotonic increasing:

$$\exp(G_k(\vec{x})) \leq \exp(|G_k(\vec{x})|) < \exp(u_g)$$

243 This result can be substituted back,

$$\left\| \vec{\nabla}_{\vec{\theta}} G(\vec{x}) \right\|_1 < \sum_{k=1}^M \frac{1}{k^2} \exp(u_g) \left\| \vec{\nabla}_{\vec{\theta}} (G_k(\vec{x})) \right\|_1$$

244 It should be noted that $G_k(\vec{x})$ and $F(\vec{x})$ have similar structure such that:

$$\left\| \vec{\nabla}_{\vec{\theta}}(G_k(\vec{x})) \right\|_1 = \left\| \vec{\nabla}_{\vec{\theta}}(F(\vec{x})) \right\|_1$$

245 This is true, even though $G_k(\vec{x}) \neq F(\vec{x})$, because the same set of basis functions are used, with
246 different coefficient parameters being the only difference. The consequence is that:

$$\left\| \vec{\nabla}_{\vec{\theta}} G(\vec{x}) \right\|_1 < \sum_{k=1}^M \frac{1}{k^2} \exp(u_g) \left\| \vec{\nabla}_{\vec{\theta}}(F(\vec{x})) \right\|_1$$

247 Substituting previously shown results gives:

$$\begin{aligned} \left\| \vec{\nabla}_{\vec{\theta}} G(\vec{x}) \right\|_1 &< \sum_{k=1}^M \frac{1}{k^2} \exp(u_g) 4nu \\ \left\| \vec{\nabla}_{\vec{\theta}} G(\vec{x}) \right\|_1 &< 4nu \exp(u_g) \sum_{k=1}^{\infty} \frac{1}{k^2} \\ \left\| \vec{\nabla}_{\vec{\theta}} G(\vec{x}) \right\|_1 &< 4nu \exp(u_g) \frac{\pi^2}{6} < 4nu\pi^2 \exp(u_g) \end{aligned}$$

248 The same argument can be used to find an upper bound for $\left\| \vec{\nabla}_{\vec{\theta}} H(\vec{x}) \right\|_1$

$$\left\| \vec{\nabla}_{\vec{\theta}} H(\vec{x}) \right\|_1 < 4nu \exp(u_h) \frac{\pi^2}{6} < 4nu\pi^2 \exp(u_h)$$

249 The original expression of interest was:

$$\begin{aligned} \left\| \vec{\nabla}_{\vec{\theta}} A(\vec{x}) \right\|_1 &\leq \left\| \vec{\nabla}_{\vec{\theta}} F(\vec{x}) \right\|_1 + \left\| \vec{\nabla}_{\vec{\theta}} G(\vec{x}) \right\|_1 + \left\| \vec{\nabla}_{\vec{\theta}} H(\vec{x}) \right\|_1 \\ \left\| \vec{\nabla}_{\vec{\theta}} A(\vec{x}) \right\|_1 &< 4nu + 4nu\pi^2 \exp(u_g) + 4nu\pi^2 \exp(u_h) \\ \left\| \vec{\nabla}_{\vec{\theta}} A(\vec{x}) \right\|_1 &< 4nu\pi^2 + 4nu\pi^2 \exp(u_g) + 4nu\pi^2 \exp(u_h) \\ \left\| \vec{\nabla}_{\vec{\theta}} A(\vec{x}) \right\|_1 &< 4nu\pi^2 (1 + \exp(u_g) + \exp(u_h)) \end{aligned}$$

250 Let the upper bound U be given by:

$$U = 4nu\pi^2 (1 + \exp(u_g) + \exp(u_h))$$

251 From the definition of the norm $\left\| \cdot \right\|_1$, and the trainable parameters θ_i with index set Θ one has that:

$$\left\| \vec{\nabla}_{\vec{\theta}} A(\vec{x}) \right\|_1 = \sum_{i \in \Theta} \left| \frac{\partial A}{\partial \theta_i}(\vec{x}) \right|$$

252 Finally,

$$\left\| \vec{\nabla}_{\vec{\theta}} A(\vec{x}) \right\|_1 = \sum_{i \in \Theta} \left| \frac{\partial A}{\partial \theta_i}(\vec{x}) \right| < U$$

253

□

254 *Remark.* For a fixed number of variables n , the model has a total of $n2^{r+2}(2M+1)$ trainable
 255 parameters. The factor of k^{-2} inside the expression for Atlas is necessary to ensure the sum is
 256 convergent in the limit of infinitely many exponential terms $M \rightarrow \infty$. Only the maximum density
 257 ($\rho = r$) cubic B-spline function has trainable parameters, so that the gradient vector is bounded in the
 258 limit of arbitrarily large densities $r \rightarrow \infty$. It is worth recalling that at most four basis functions are
 259 active for uniform cubic B-spline functions, regardless of the density, but the smaller densities cannot
 260 be trainable, otherwise this property does not hold. The gradient vector has bounded norm for any
 261 number of basis functions and exponential terms. The bounded gradient vector implies that Atlas is
 262 numerically stable during training, regardless of its size or parameter count.

263 B.5.4 Atlas distal orthogonality

264 **Property 3** (Distal orthogonality). For any Atlas model $A(\vec{x})$ and $\forall \vec{x}, \vec{y} \in D(A) \subset R^n$ and
 265 trainable parameters θ_i , there exists a $\delta > 0$ such that:

$$\min_{j=1, \dots, n} \{|x_j - y_j|\} > \delta \implies \langle \vec{\nabla}_{\vec{\theta}} A(\vec{x}), \vec{\nabla}_{\vec{\theta}} A(\vec{y}) \rangle = 0$$

266 *Proof.* Let $A(\vec{x})$ denote some Atlas model, with mixed-density B-spline functions $f_j(x_j)$, $g_{i,j}(x_j)$,
 267 and $h_{i,j}(x_j)$ in the form:

$$\begin{aligned} A(\vec{x}) &= \sum_{j=1}^n f_j(x_j) + \sum_{k=1}^M \frac{1}{k^2} \exp(\sum_{j=1}^n g_{k,j}(x_j)) - \frac{1}{k^2} \exp(\sum_{j=1}^n h_{k,j}(x_j)) \\ &= F(\vec{x}) + \sum_{k=1}^M \frac{1}{k^2} \exp(G_k(\vec{x})) - \frac{1}{k^2} \exp(H_k(\vec{x})) \\ &= F(\vec{x}) + G(\vec{x}) - H(\vec{x}) \end{aligned}$$

268 With each mixed-density B-spline function in $A(\vec{x})$ given by:

$$\begin{aligned} f_j(x_j) &= \sum_{\rho=0}^r \sum_{i=1}^{2^{\rho+2}} \theta_{f,(\rho,i,j)} S_{\rho,i}(x_j) \\ g_{k,j}(x_j) &= \sum_{\rho=0}^r \sum_{i=1}^{2^{\rho+2}} \theta_{g,(\rho,i,k,j)} S_{\rho,i}(x_j) \\ h_{k,j}(x_j) &= \sum_{\rho=0}^r \sum_{i=1}^{2^{\rho+2}} \theta_{h,(\rho,i,k,j)} S_{\rho,i}(x_j) \end{aligned}$$

269 Any mixed-density functions Φ and Ψ that act on different components of the input must have
 270 orthogonal parameter gradients, since each input variable has its own associated parameters:

$$\langle \vec{\nabla}_{\vec{\theta}} \Phi(x_i), \vec{\nabla}_{\vec{\theta}} \Psi(y_j) \rangle = 0 \quad \forall i \neq j$$

271 Generally, since all mixed-density functions have parameters that are independent of each other it
 272 follows that for any mixed-density B-splines Φ and Ψ :

$$\langle \vec{\nabla}_{\vec{\theta}} \Phi(x_j), \vec{\nabla}_{\vec{\theta}} \Psi(y_j) \rangle = 0 \quad \forall \Phi \neq \Psi$$

273 Thus, one need only compare the parameter gradients of each mixed-density B-spline function Ψ
 274 with itself. The inner-product of the parameter gradient of Ψ evaluated on two different inputs is
 275 given by:

$$\langle \vec{\nabla}_{\vec{\theta}} \Psi(x_j), \vec{\nabla}_{\vec{\theta}} \Psi(y_j) \rangle$$

276 The inner-product given above is not zero in general. However, as illustrated in Figure 18, for any
 277 mixed-density B-spline function Ψ there exist a $\delta > 0$, such that:

$$|x_j - y_j| > \delta \implies \langle \vec{\nabla}_{\vec{\theta}} \Psi(x_j), \vec{\nabla}_{\vec{\theta}} \Psi(y_j) \rangle = 0$$

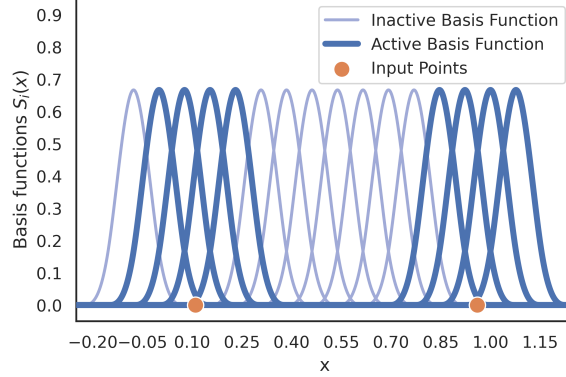


Figure 18: Visual proof of distal orthogonality for single-variable ρ -density B-splines.

278 This is because each basis function is zero everywhere, except on some small sub-interval. If this is
 279 true for all $j = 1, \dots, n$, then the parameter gradients evaluated at \vec{x} and \vec{y} must be orthogonal. If this
 280 is true for all $j = 1, \dots, n$, then it is true for the minimum. The converse is true by transitivity such
 281 that:

$$|x_j - y_j| > \delta \forall j = 1, \dots, n \iff \min_{j=1, \dots, n} \{|x_j - y_j|\} > \delta$$

282 Finally,

$$\min_{j=1, \dots, n} \{|x_j - y_j|\} > \delta \implies \langle \vec{\nabla}_{\vec{\theta}} A(\vec{x}), \vec{\nabla}_{\vec{\theta}} A(\vec{y}) \rangle = 0$$

283

□

284 *Remark.* Two points that sufficiently differ in each input variable have orthogonal parameter gradients.
 285 It is worth mentioning that the condition resembles a cross-like region in two variables, and planes that
 286 intersect in higher dimensions. Distal orthogonality means Atlas is reasonably robust to catastrophic
 287 forgetting.