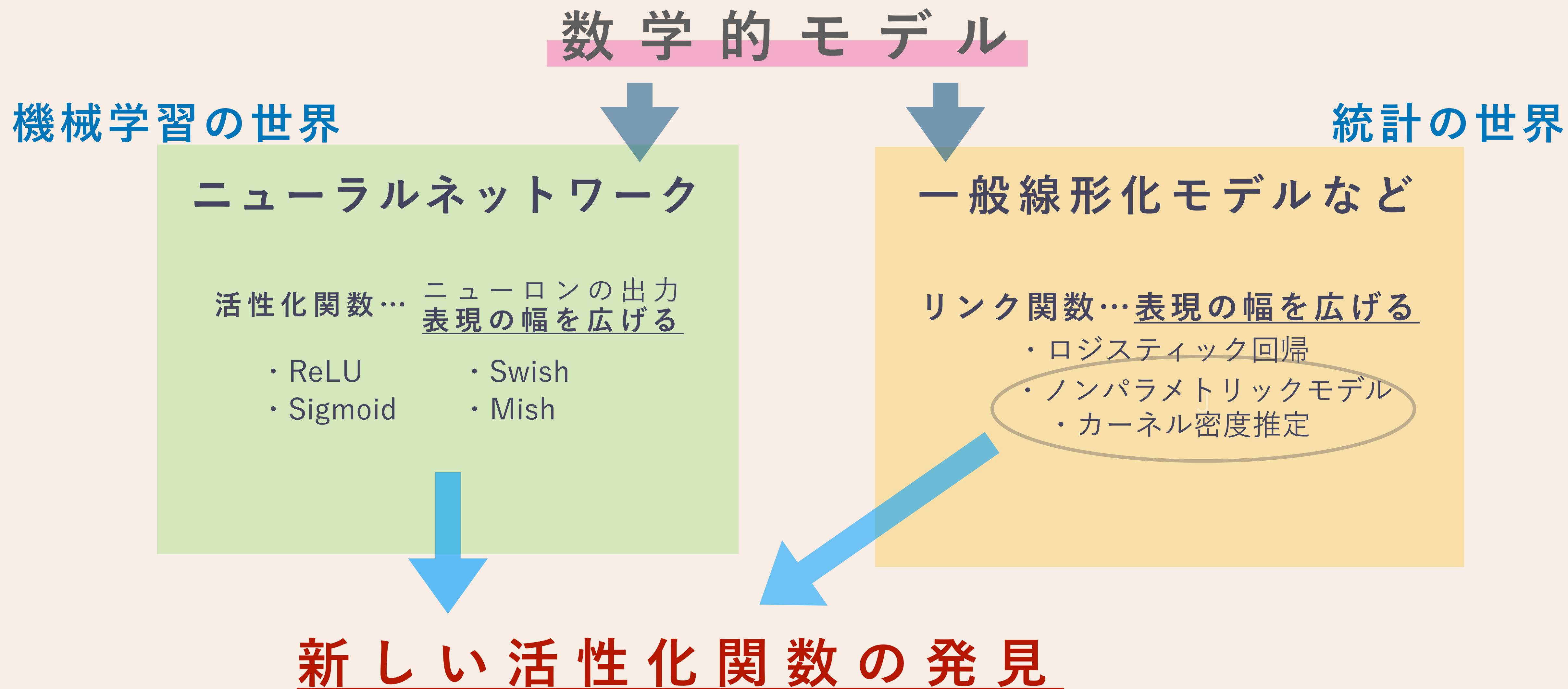


カーネル密度推定を用いた汎用的な活性化関数

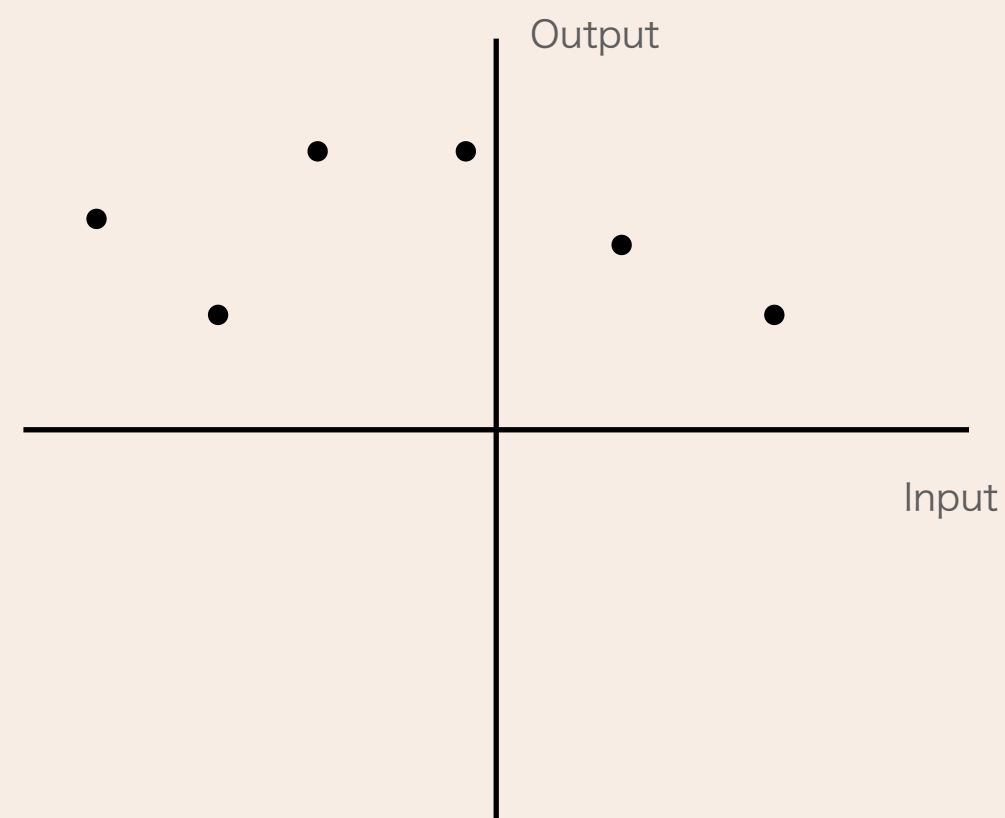
neco B4 vapor
vapor@sfc.wide.ad.jp
Adviser ks91



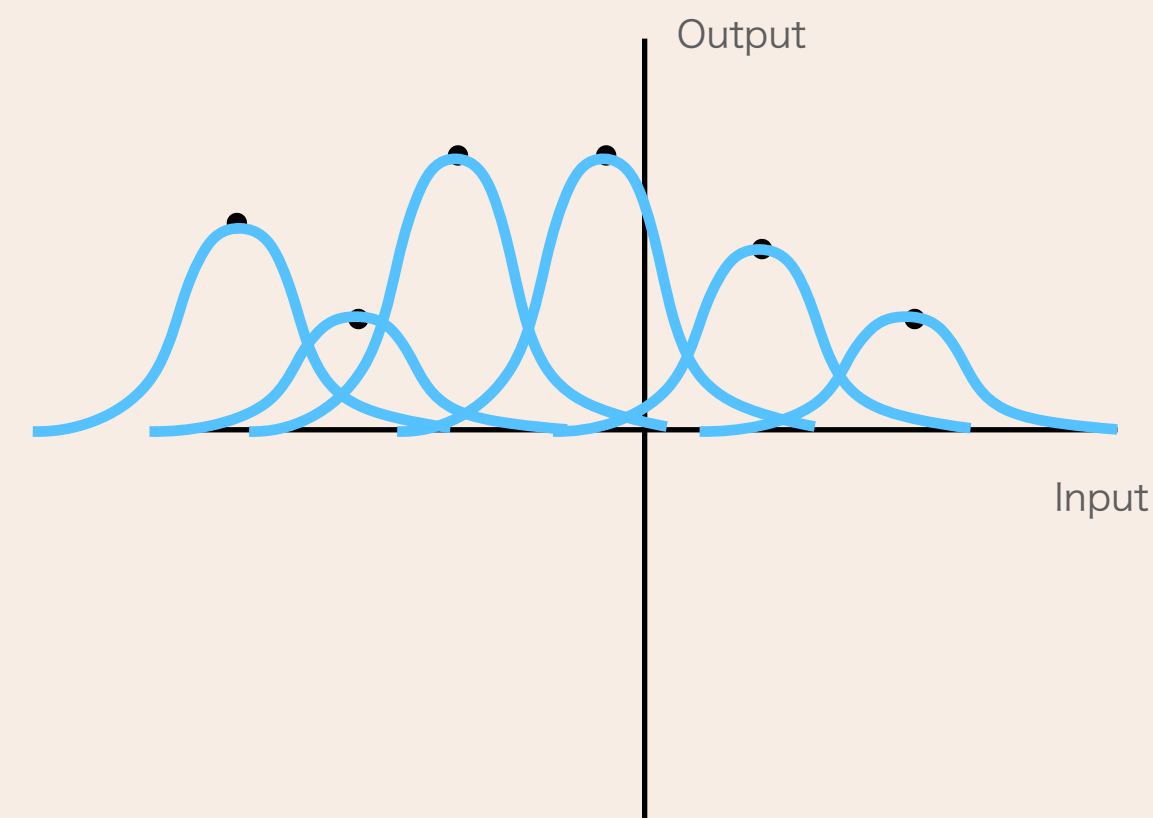
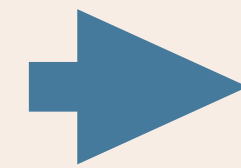
背景. 1 (リンク関数のノンパラメトリック推定)

カーネル密度推定

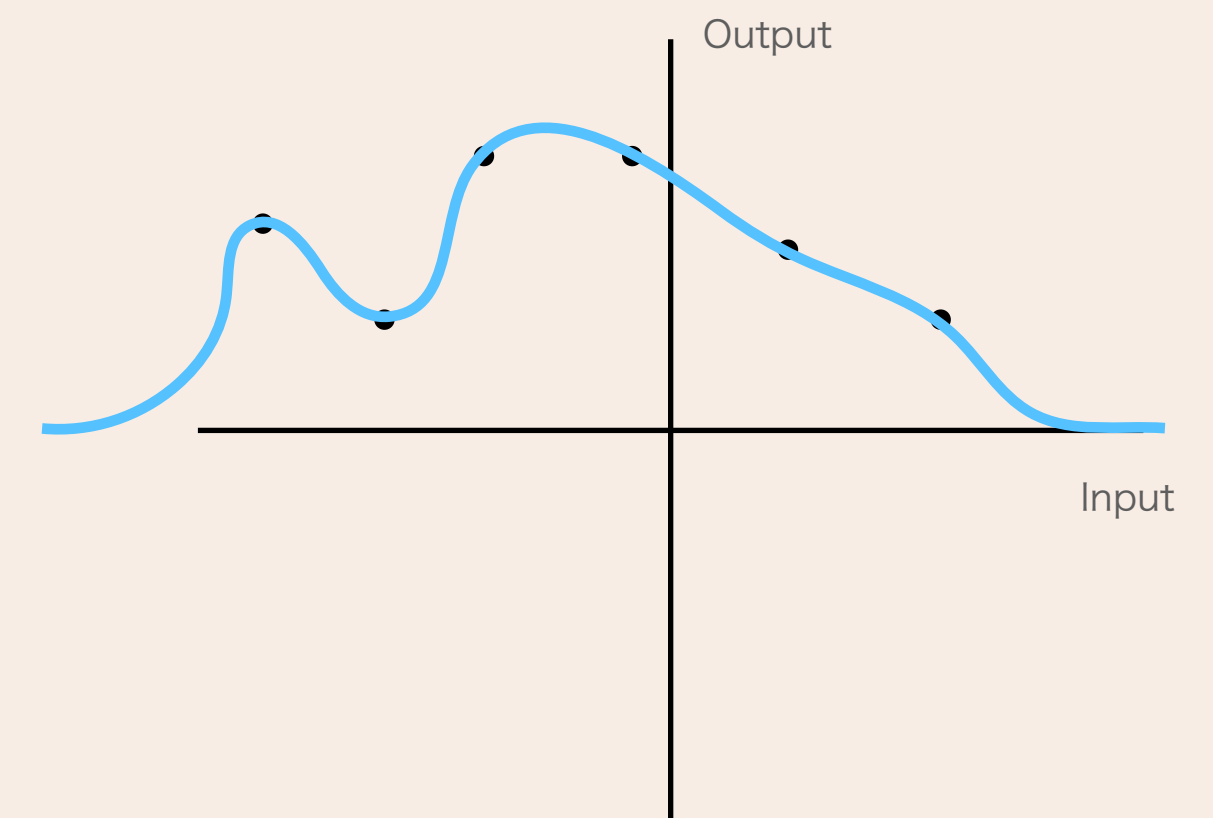
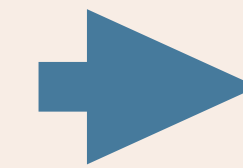
- ・ ノンパラメトリックモデルの一つ
- ・ データ点の周囲にカーネル関数を置き、関数の近似を行う



データ点がまばらに存在する。



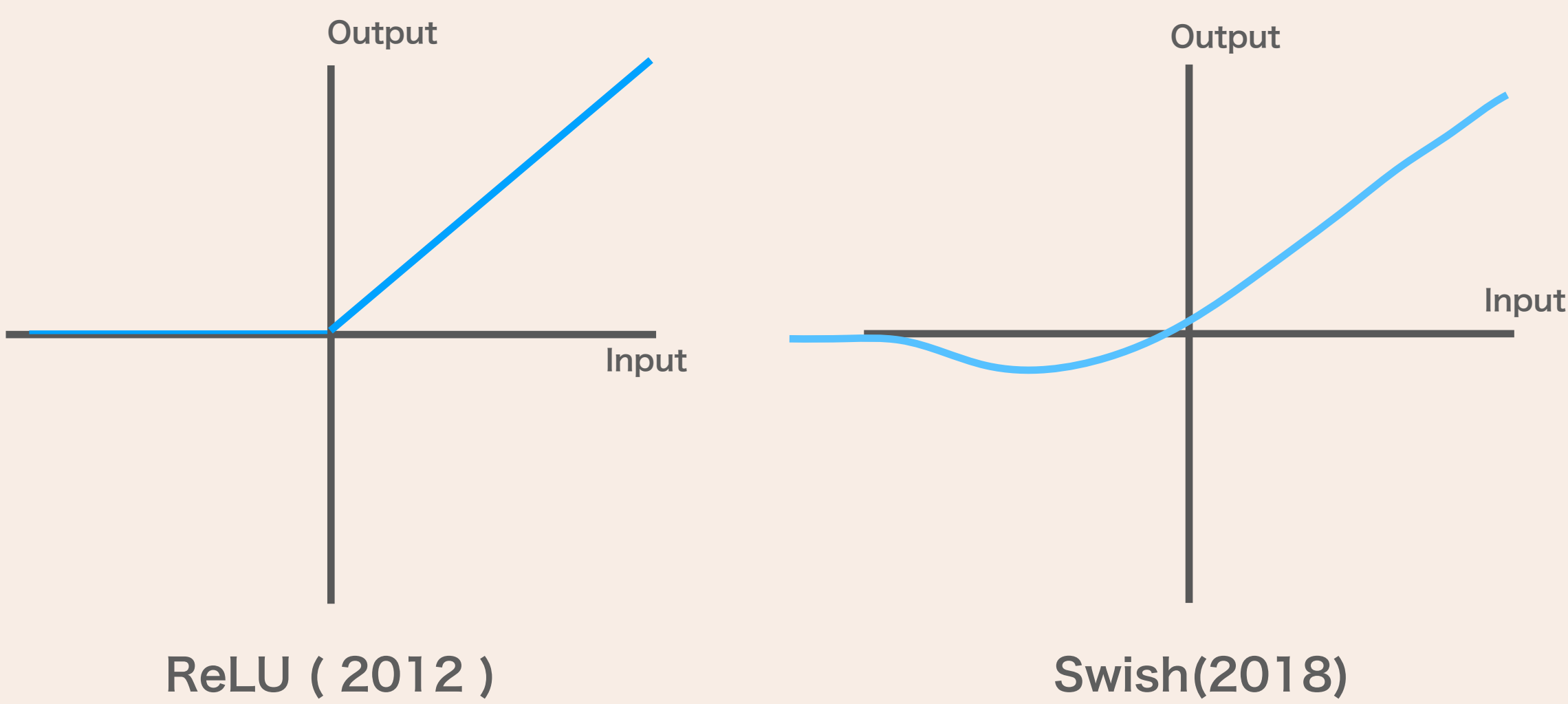
周囲にカーネル関数を置く



カーネル関数を足し合わせて関数の近似を行う

背景（活性化関数の問題点）

精度の向上問題

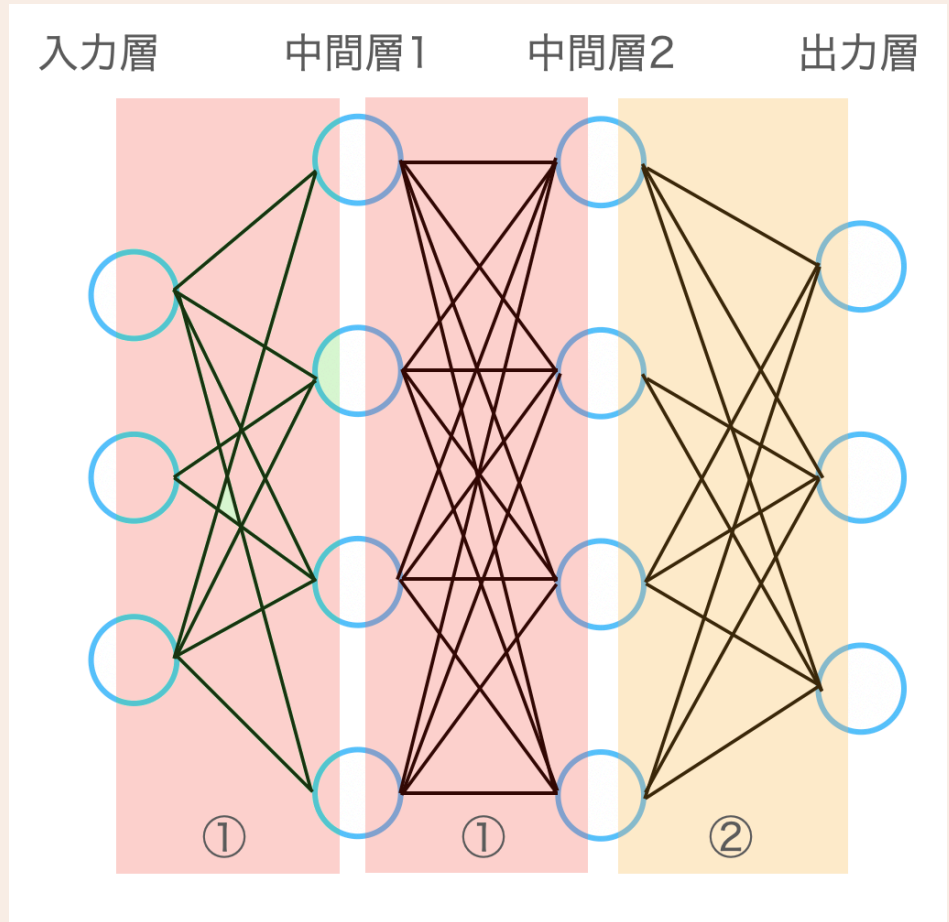


関数の形状が単調増加のシンプルなものから、複雑な形状の方が精度が上がる



より精度の高い活性化関数が求め続けられている

活性化関数の選択問題



	①	②
分類	ReLUなど	ReLUなど
回帰	ReLUなど	Sigmoidなど

「分類」や「回帰」といったデータの出力形式によって、出力層に使われる活性化関数が変わる



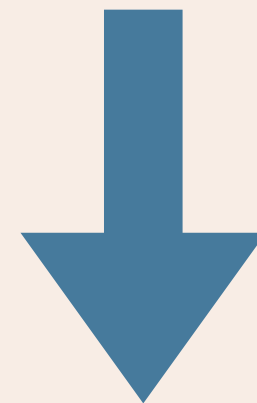
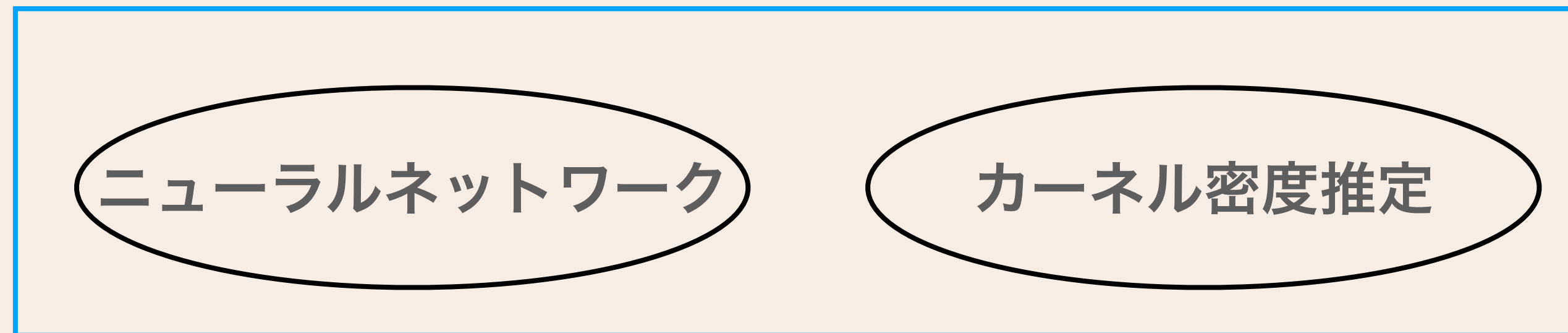
関数の選択は、経験則に基づく必要がある、何がいい選択かわからない

本研究が取り組むべき課題

1. 高い汎用度で活性化関数が自動で推論でき、既存のものより安定的で精度が高くなること
2. ノンパラメトリックに活性化関数が推定でき、既存の扱っていた活性化関数との差分を確認できる ようにすること

提案手法

統計学の知識を組み合わせながら、ニューラルネットワークに使える新しい活性化関数を考える



新しいアルゴリズムの提唱
K-AF (Kernel Activation Function)

K-AFの式

Ichimura(1993)の提案を変形する[1]

$$\tilde{G}(\mathbf{X}_i, \mathbf{W}, \mathbf{X}^{calc}, \mathbf{Y}^{calc}) = \frac{\sum_{i \neq j} K\left(\frac{\mathbf{X}_j^{calc} \mathbf{W} - \mathbf{X}_i \mathbf{W}}{h_{calc}}\right) \mathbf{Y}_j^{calc}}{\sum_{i \neq j} K\left(\frac{\mathbf{X}_j^{calc} \mathbf{W} - \mathbf{X}_i \mathbf{W}}{h_{calc}}\right)}$$

\mathbf{X}^{calc} ...データセットの一部を表し、 $\mathbf{X}^{calc} \subset \mathcal{D}_x$
 \mathbf{Y}^{calc} ...データセットの一部を表し、 $\mathbf{Y}^{calc} \subset \mathcal{D}_y$

K ...カーネル関数（本実験ではガウス関数）
 \mathbf{W} ...重み
 h_{calc} ...バンド幅

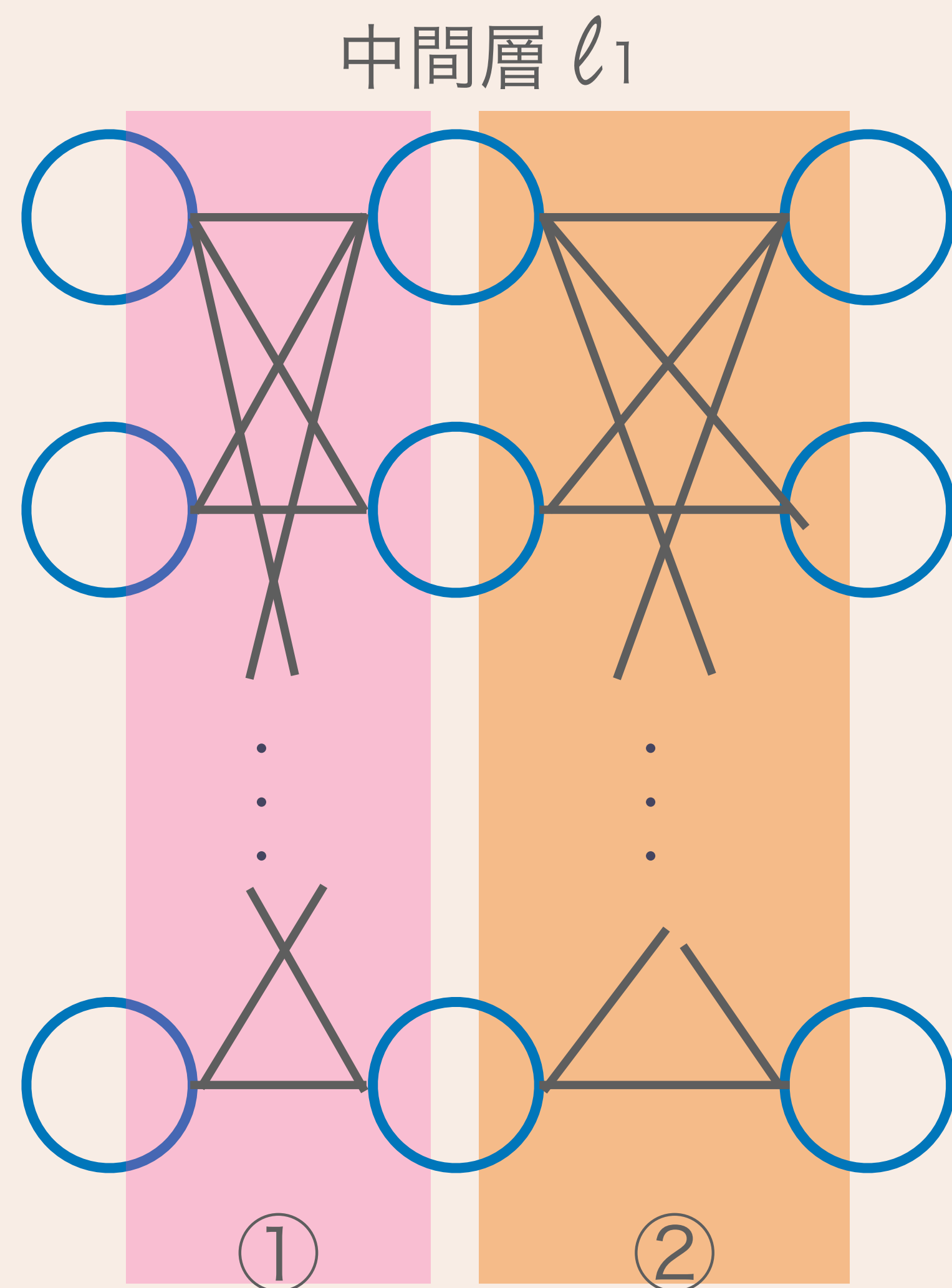
実験項目

- 1. K-AFの精度の評価
- 2. K-AFの関数の可視化による、既存の関数の妥当性の評価
- 3. K-AFが勾配発散しないための条件調査

本研究を動かすための実行環境及びライブラリ

環境	バージョン
Python	3.8.5
scikit-learn	0.23.2
numpy	1.18.5
pytorch	2.5.0

項目	仕様
CPU	2.6GHz Intel core i7
メモリ	16GB 2400 MHz DDR4
ストレージ	Macintosh HD 256GB
OS	macOS High Sierra 10.13.6



今回のニューラルネットワークは
中間層を一つで実験し、出力層にK-AFを配置し
実験1~3を評価する

- ・ 図2の①にはReLUを置く
- ・ ②にK-AF等の様々な活性化関数を配置する

図2:本実験で使うニューラルネットワークの構成

実験. 1の概要

K-AFの精度の評価

- 1. 図2の②に対して
ReLU, Sigmoid, Tanh, Swish, Mish, K-AF
の6つの活性化関数を入れ替え、K-AFの評価をする
- 2. 比較用のデータの性質及び設定を以下の表に示す

データセット名	サンプル数	入力の次元	出力の次元	出力の形式	中間層の数 l_1
✕ iris	150	3	3	分類	4
digits	1797	3	64	分類	100
✕ wine	178	3	13	分類	40
✕ boston	506	13	1	回帰	20
breast_cancer	569	30	2	分類	30

※スライドではdigitsとbostonの性能を確認する

実験. 1.1 digitsの比較設定

以下の表の二つの設定でdigitsを学習した際のK-AFのAccuracy及びValidationLossを取得する

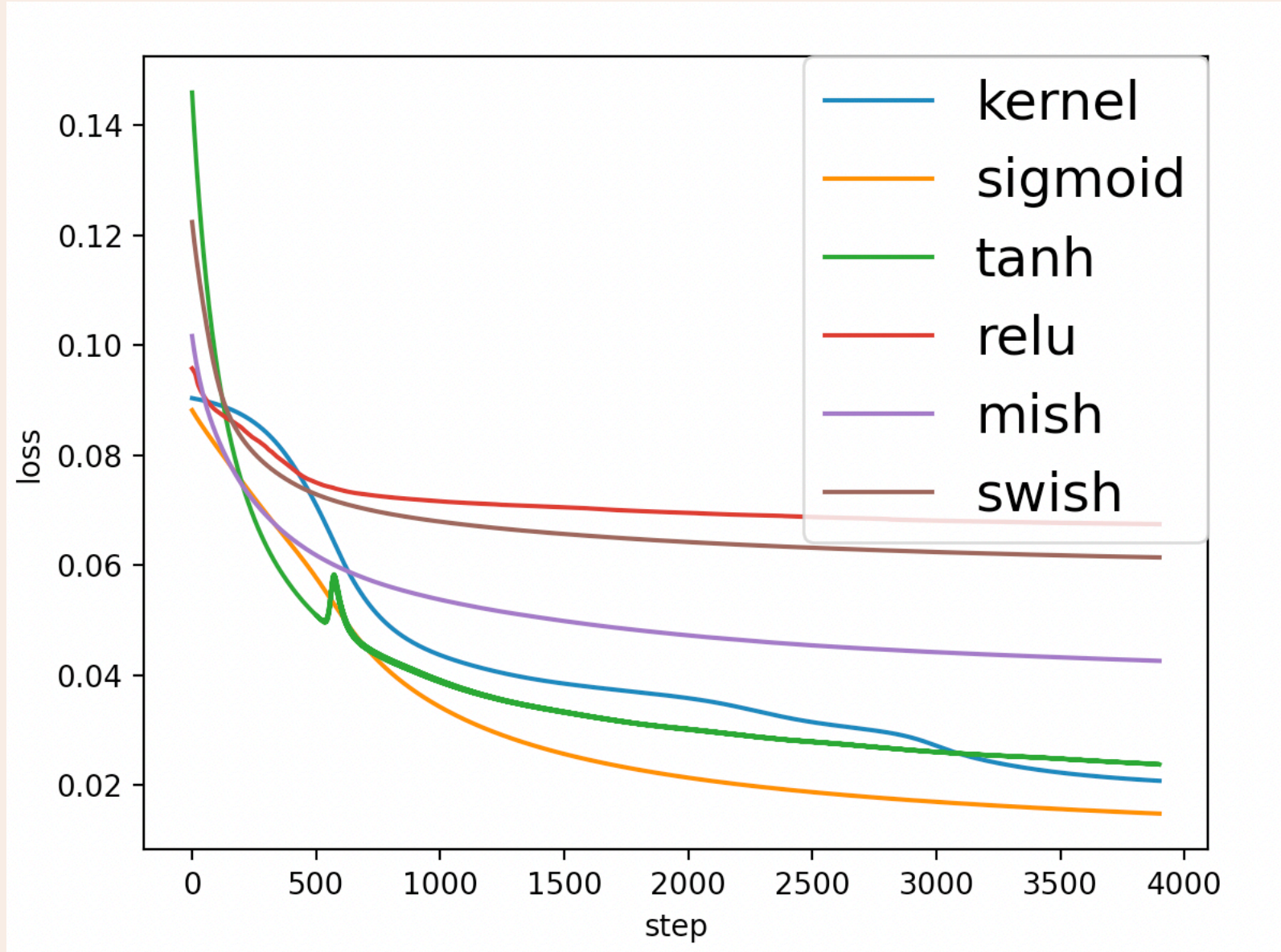
設定名	設定 1	設定 2
LearningRate	10^{-2}	10^{-3}
Initializer	KaimingUniform	KaimingUniform
Optimizer	SGD	SGD
Regularizer	なし	なし
epoch_num	2000	10000
exp_num	3	3
calc_num	36	36

Bostonのデータセットで比較実験をする際の設定

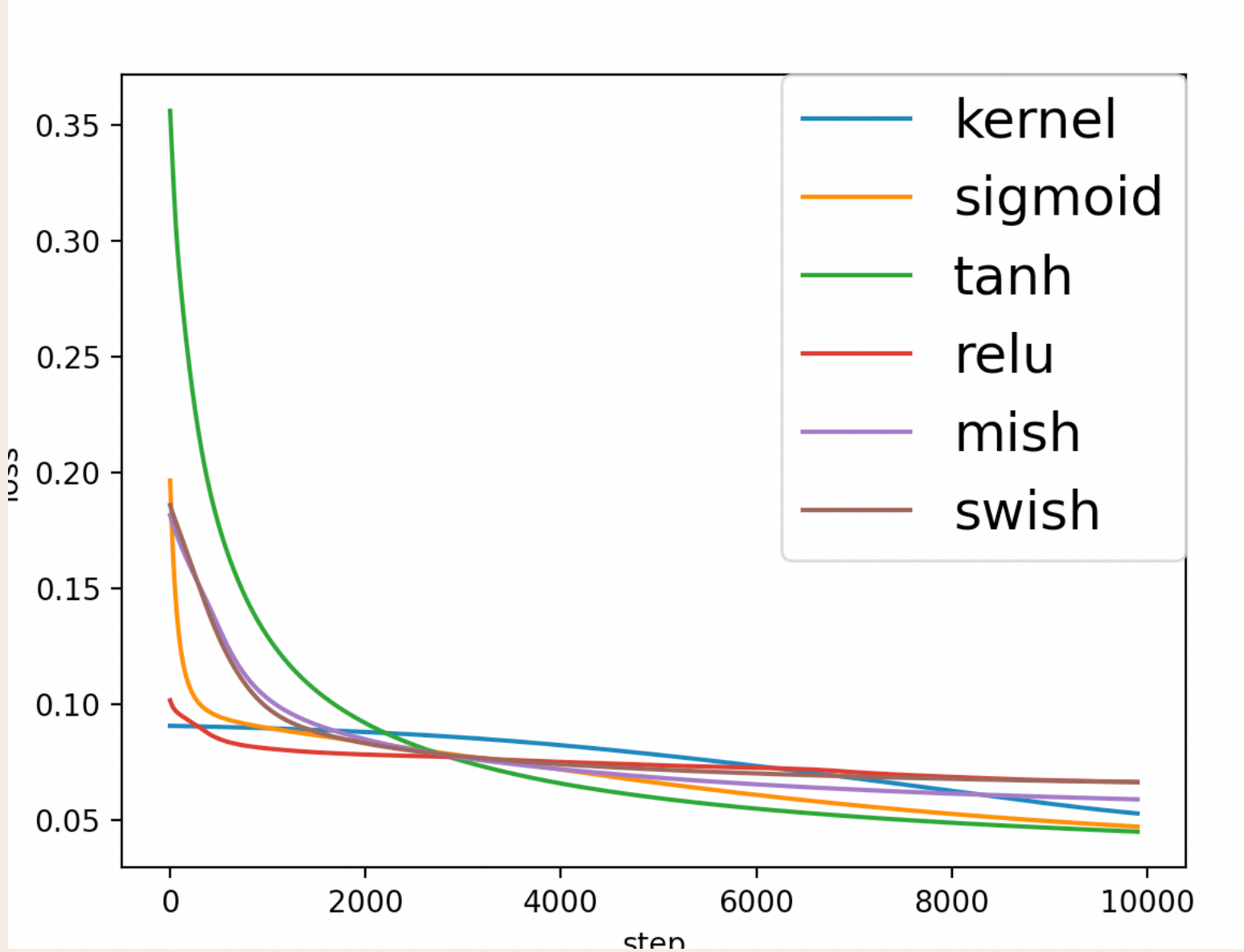
※Accuracyはexp_numの数の平均値を取得する

epoch_numはepoch数、exp_numは実験回数、calc_numは \mathbf{X}^{calc} の数

実験. 1.1 digitsの比較結果



設定1のValidationLoss



設定2のValidationLoss

活性化関数	設定 1 の Accuracy	設定 2 の Accuracy
K-AF	91.6	60.0
Sigmoid	92.0	68.6
Tanh	95.6	89.6
ReLU	38.0	54.3
Swish	50.3	56.6
Mish	55.6	53.3

各活性化関数のAccuracyによる性能比較。他の関数と同程度の性能を出した。

1. AccuracyもLossも既存のSigmoidやTanhなどと同程度の性能
2. 特に学習率が高い場合は高い性能を出している

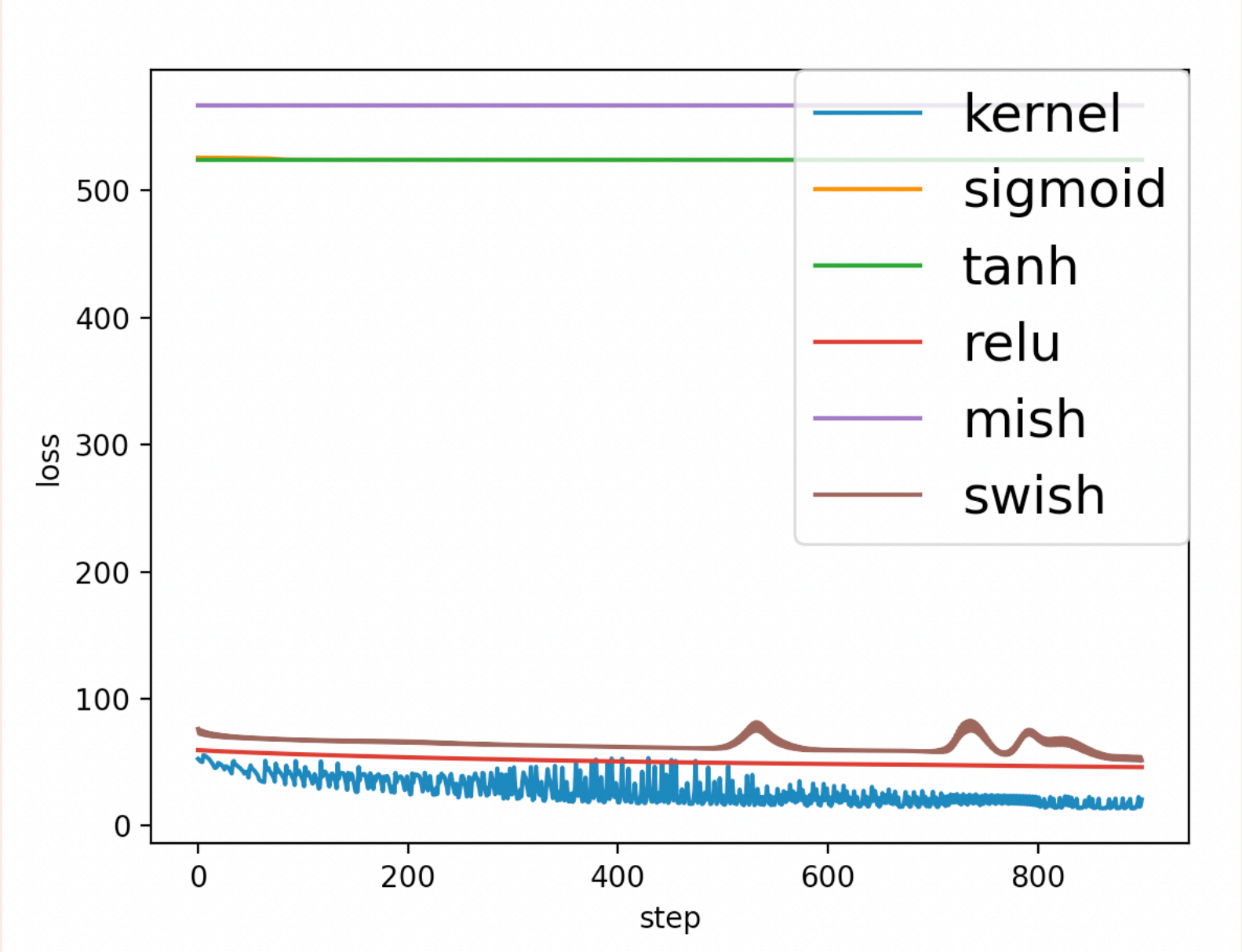
実験. 1.2 bostonの比較設定

以下の表の二つの設定でBostonを学習した際のK-AFのMSE及びValidationLossを取得する。

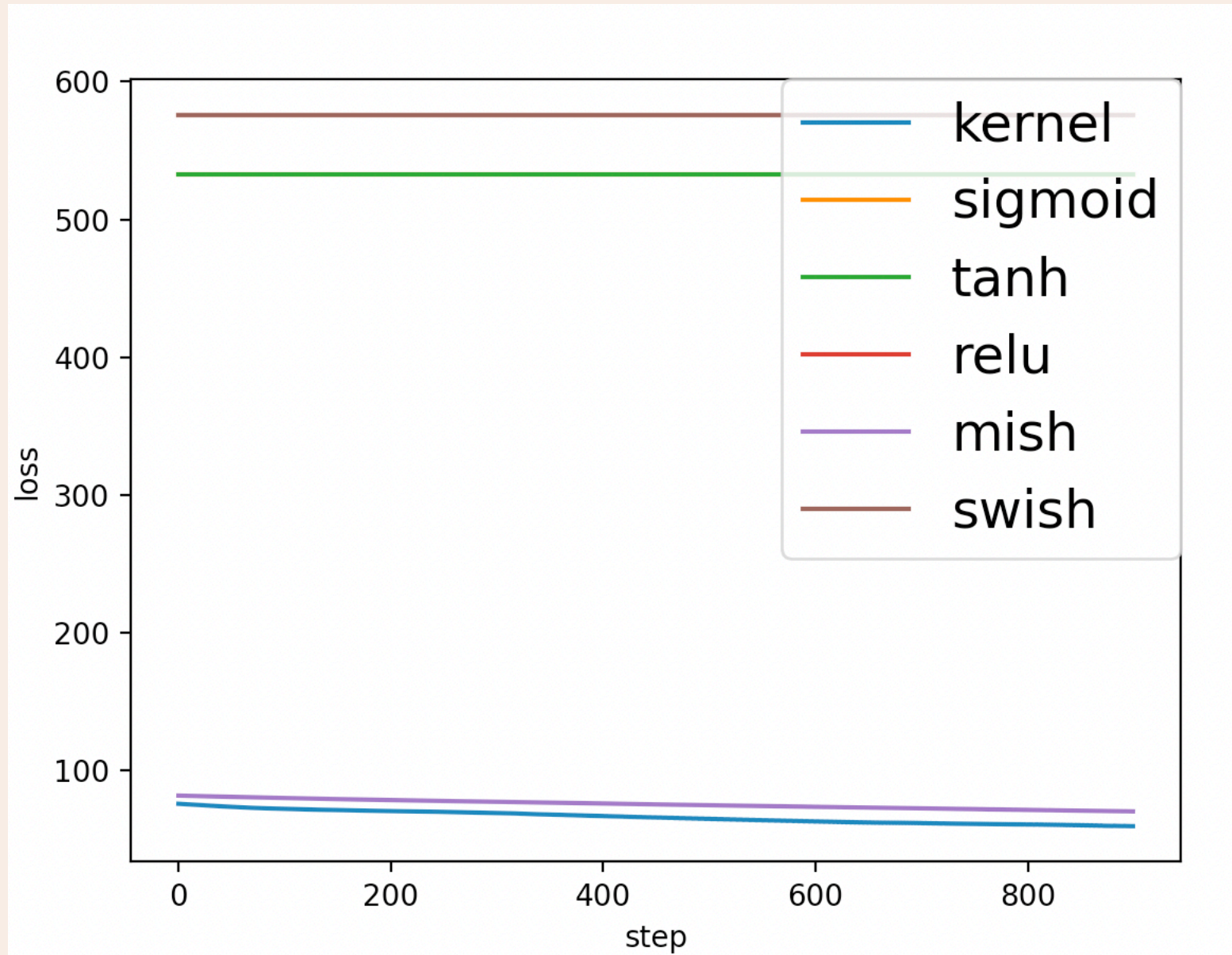
設定名	設定 1	設定 2
LearningRate	10^{-5}	10^{-5}
Initializer	KaimingUniform	Xavier
Optimizer	SGD	SGD
Regularizer	なし	なし
epoch_num	1000	1000
exp_num	5	5
calc_num	26	26

Bostonのデータセットで比較実験をする際の設定
※Bostonは回帰問題なのでMSE(最小二乗誤差)で性能評価を行う事に注意

実験. 1.2 bostonの比較結果



設定1のValidationLoss



設定2のValidationLoss

活性化関数	設定 1 の MSE	設定 2 の MSE
K-AF	49.4	69.4
Sigmoid	523.0	541.4
Tanh	540.6	567.8
ReLU	359.2	473.4
Swish	257.0	372.4
Mish	360.0	472.4

各活性化関数のMSEによる性能比較
K-AFが優れている。

1. Bostonは正確に推測するためには決定木など十分な複雑性を持ったモデルが必要なため、K-AFでは非常に高い性能を出した
2. MishやSwishでも回帰問題のため、高い性能を出している

実験. 2の概要

K-AF関数の可視化による、既存の関数の妥当性の評価

- 1. 推論した活性化関数の形状を調査し、既存の活性化関数との違いを比べる
- 2. 定性評価を行うために、以下の表項目を用い調査をする

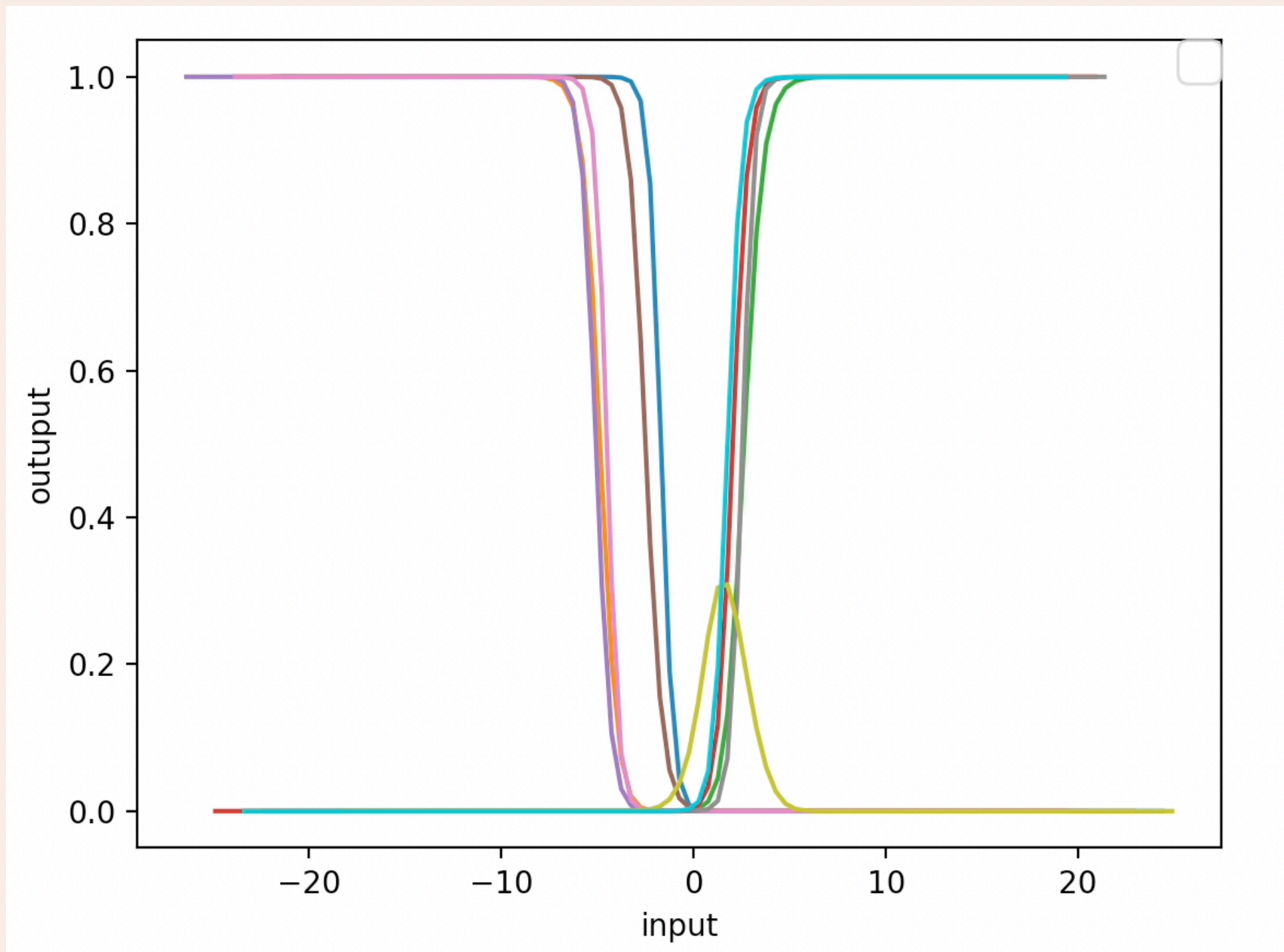
単調増加関数か	上限値があるか
○ or ×	○ or ×

- 3. 各データセットにおいて以下の表項目の設定で学習を行う

データセット名	LearningRate	Optimizer	Initializer	Reguralizer	中間層の数	calc_num	epoch_num
✕ iris	10^{-1}	SGD	KaimingUniform	なし	4	30	1000
digits	10^{-1}	SGD	KaimingUniform	なし	100	36	1000
✕ wine	10^{-2}	SGD	KaimingUniform	なし	40	36	1000
boston	10^{-5}	SGD	KaimingUniform	なし	20	26	1000
breast_cancer	10^{-2}	SGD	KaimingUniform	なし	30	285	1000

※実験1同様にスライドではdigitsとbostonの場合のみ観測する

実験. 2.1 digitsの場合(画像)



digitsの学習で生成されたK-AF

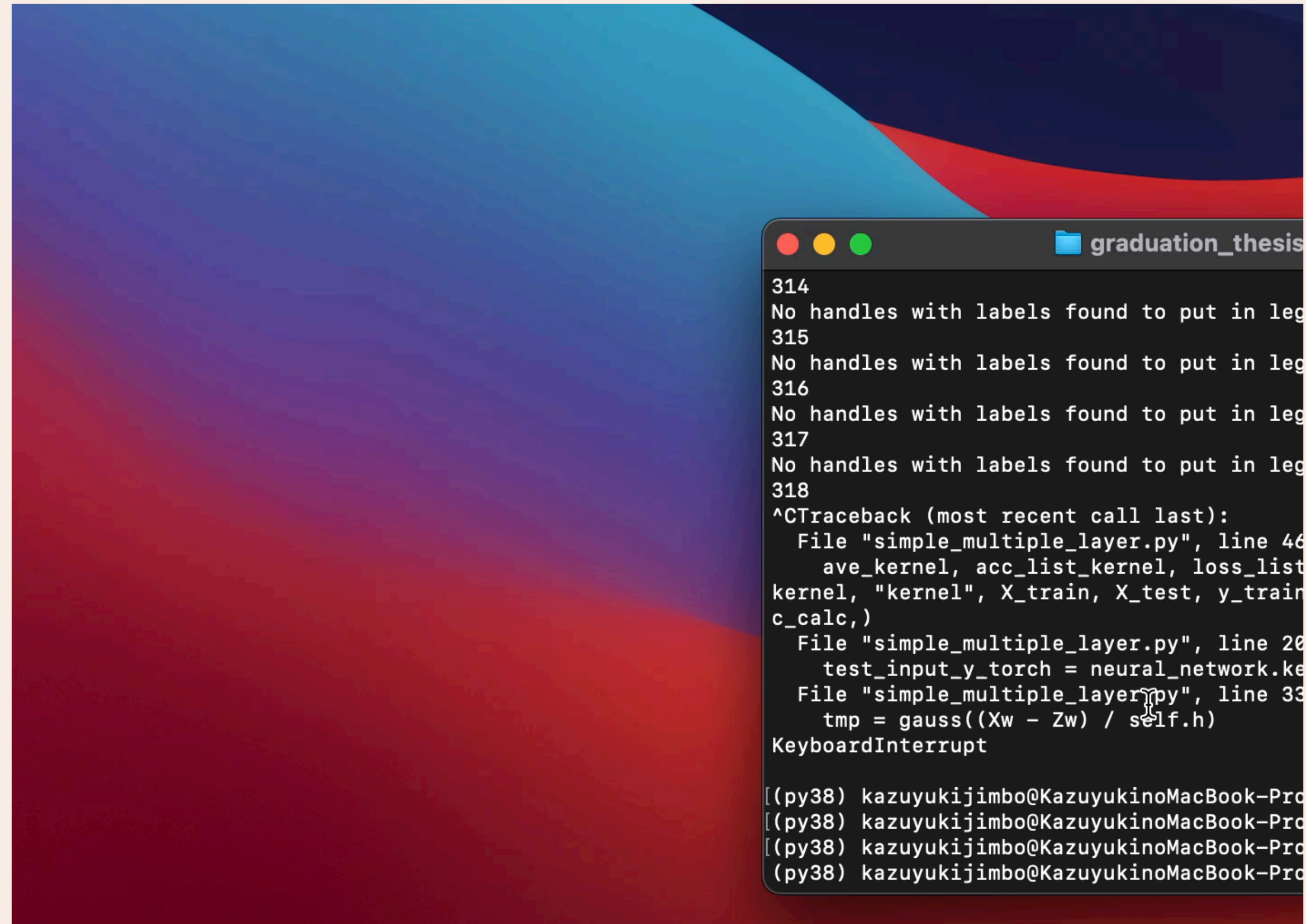
digitsの学習で生成されたK-AFの形状評価

単調増加関数か	上限値があるか
▲	○

1. ほぼ全てのK-AF 単調増加関数もしくは単調減少関数となった
2. ラベリング問題なので上限値が1.0になるのも想定通りである
3. Sigmoidに等に性質や形が似ているので、このデータセットではSigmoidでも十分な性能持っているとも言える

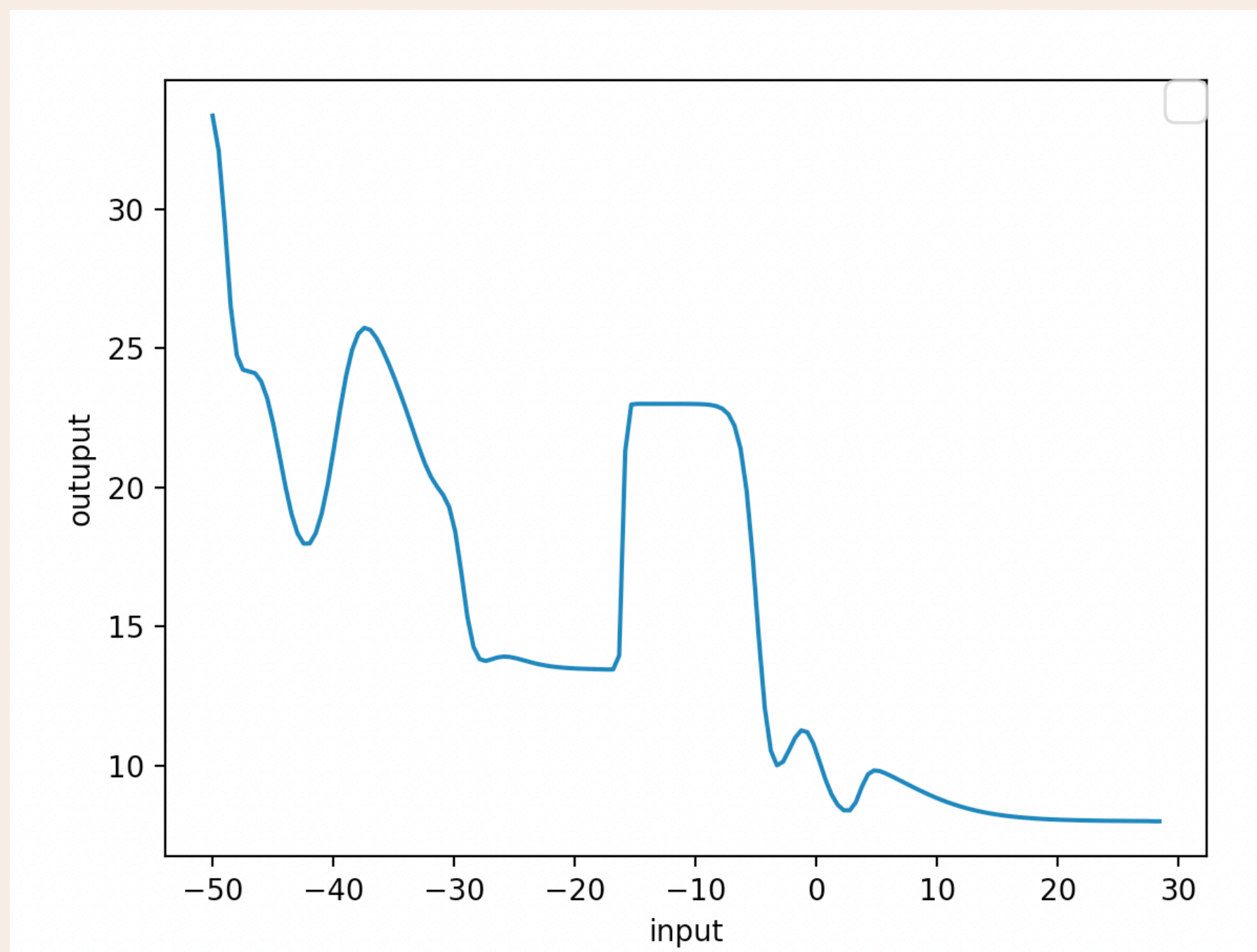
実験. 2.1 digitsの場合(動画)

digitsのK-AFの学習過程の可視化



<https://github.com/latte0/RG-Thesis-Template/blob/main/asset/digits.mov?raw=true>

実験. 2.2 bostonの場合(画像)



bostonの学習で生成されたK-AF

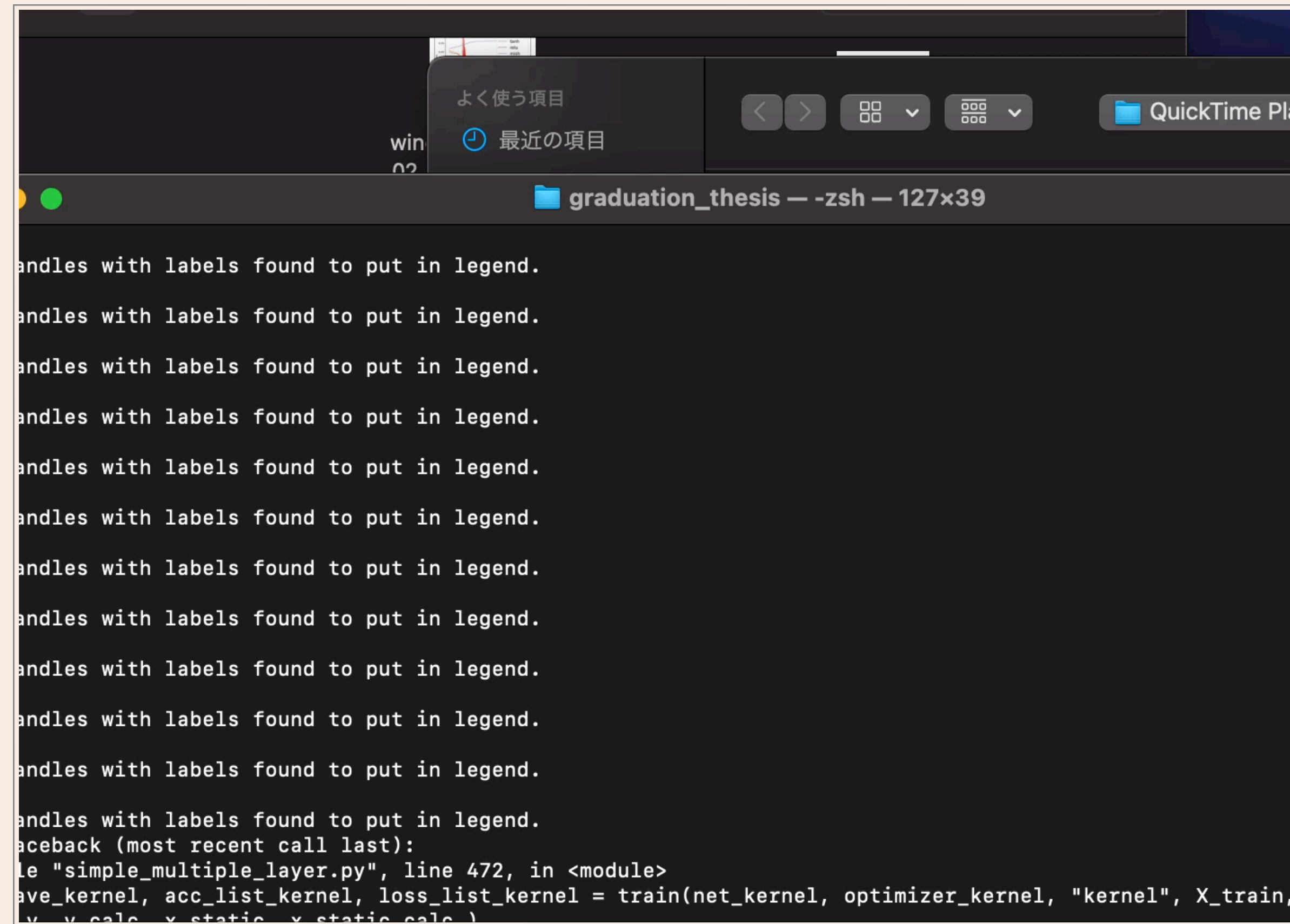
bostonの学習で生成されたK-AFの形状評価

単調増加関数か	上限値があるか
×	▲

1. 歪な形をしていて、単調増加なReLUやSigmoidでは表現のとして不十分である
2. 回帰問題であるので、上限値がないような関数が適切であった

実験. 2.2 bostonの場合(動画)

BostonのK-AFの学習過程の可視化



<https://github.com/latte0/RG-Thesis-Template/blob/main/asset/boston.mov?raw=true>

実験. 3の概要

K-AFが勾配発散しないための条件調査

- 1. 以下の表の変更リストの全ての組み合わせで勾配爆発がどの程度起こるかテストをする。

ニューラルネットワークの構成	変更のリスト
Initializer	[Xavier, KaimingUniform]
Regularizer	[何もしない, L1 ノルム, L2 ノルム]
Optimizer	[SGD, RMSProp, AdaGrad, Adam]

- 2. データセットは**boston**を用い、その他のパラメータを以下に設定する。

ハイパーパラメータ	値
LearningRate	10^{-5}
中間層の次元数 l_1	20
epoch_num	5
exp_num	1000
calc_num	50

※論文中のクリッピングの実験は特に効果がなかったのでスライドでは省く。

実験. 3の結果

全てのパターンの中でも最も勾配爆発確率が低かったトップ5

順位	Initializer	Optimizer	Reguralizer	勾配爆発回数	勾配爆発確率
1	Xavier	Adam	non	275	27.5%
2	Xavier	Adam	l1	290	29.0%
3	Xavier	Adam	l2	294	29.4%
4	Xavier	SGD	l1	294	29.4%
5	Xavier	RMSprop	l1	298	29.8%

1. XavierはKaimingUniformの2倍程度、勾配爆発の確率が低かった
2. OptimzerとしてAdamが上位を占めた

実験結果まとめ

1. 一部の条件で、**回帰と分類の双方で高い精度**の活性化関数をK-AFで導くことができた
2. K-AFの形状の調査により、**既存の活性化関数は問題に対して適切とは言えない**場合があることが判明した
3. K-AFは**初期値に依存して勾配爆発**が防げることがわかった

解決した課題

1. 出力層に用いる活性化関数を状況に応じた適切な形に変わる汎用的な関数を導き、高い精度を出せるようにした
2. 各問題やデータに応じた最適な活性化関数の形が模索できるようになった
3. 分類や回帰といった問題を考えずとも、高い精度でニューラルネットワークが構築できるようになった

本研究の問題点と展望

勾配の発散問題

データセットが複雑な場合、勾配が発散するのが一定の確率で防げてない

計算時間の問題

使用するデータ数が多いと計算が重くなってしまう

ディープラーニングへの応用

出力層のみでなく、中間層の活性化関数も推論できる関数を見つける必要がある

- [1] Hidehiko Ichimura, Wolfgang Hardle, and Peter Hall. Optimal smoothing in single index model. J. Econometrics, 58(1-2):71–120, 1993.
- [2] Adam Tauman Kalai and Ravi Sastry. The isotron algorithm: High-dimensional isotonic regression. In COLT '09, 2009.
- [3] Xavier. Glorot and Antoine. Bordes. Deep sparse rectifier neural networks. 14th International Conference on Artificial Intelligence and Statistics, 2011.
- [4] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. arXiv preprint arXiv:1512.03385, 2015.
- [5] Pytorch. <https://pytorch.org/>, 2021. (Accessed on 02/03/2021)
- [6] scikit-learn Machine Learning in Python, <https://scikit-learn.org/stable/> (Accessed on 02/03/2021)