

Random Fourier Features

石川 徹也 <tiskw111@gmail.com>

2020 年 3 月 3 日

はじめに

本文書は Rahimi *et al.* [1] によって 2007 年に提案された Random Fourier Features (以下, RFF と略す) について, カーネル法をある程度ご存じの読者を対象に解説した文書です. RFF は, 端的に言えばカーネル法におけるカーネル関数を有限次元近似する手法であり, カーネル法をより高速に, そして大規模なデータセットに適用することができるようになる手法です.

本文書では, 記号の確認も含めてカーネル法について簡単にご説明はしていますが, 初めてカーネル法に触れる読者には不十分極まりないと思います. カーネル法をご存じない方は, 事前に赤穂 [3] や Bishop [4] などに目を通して頂くことをお勧めします.

また, 本文書は意欲的な高校生や大学生諸君にとっても有用だと思います. 大学初年度の教養数学をきちんと修めていれば, 本文書を理解することはさほど困難ではありません. 少し背伸びをすれば高校生諸君でも十分に理解できる内容です. さらに RFF は, 2007 年に提案され, すでに 2000 件近くの引用がなされ, 現在でも改良が続けられている手法です. すでに 10 年以上昔の成果ですので, 変化の激しい機械学習の分野ではもはや新しい結果とは言えません. それでも学生の皆さんが普段勉強している内容と比較すれば極めて最近の, そして現在でも研究対象として注目され続けている結果です. 大昔の結果ばかり勉強させられ, 世界の最先端がどこにあるのか分からなくなってしまった方々にとって, この文書は良いきっかけを提供できるかもしれません.

また, 本文書の最後では, 深層学習との関係性についても簡単に触れています. というのも, 本文書を書くきっかけは

「カーネルサポートベクターマシンを深層学習のフレームワークで実装したいんだが, どうすりゃ良いんだコノヤロー！」

という質問(挑発?)を職場の同僚から受けたことに端を発しているためです. この問いもまた RFF を用いることで肯定的に解決できます.

何はともあれ, まずは RFF の目的と意義から順に見ていくことにしましょう.

NOTE: 同僚の名誉のために述べておくと, 上述の深層学習のフレームワークとは Tensorflow や PyTorch のような有名なワークステーション向けフレームワークではなく, 組み込み向けのマイナーなものです. このフレームワークは, ニューラルネットワークの実装を主目的に設計されて

いること, 組み込み向けゆえにサポートされている機能が極めて限定的なことから, カーネルサポートベクターマシンのようなニューラルネットワークでない機能を実装するのは容易ではないという事情があります. Tensorflow や PyTorch のような柔軟性の高いフレームワークであれば RFF を用いずとも実現は簡単です.

1 なぜ RFF が必要なのか?

機械学習におけるクラス分類問題を解く手法のひとつとして, カーネルサポートベクターマシンと呼ばれる手法がある. これは通常のサポートベクターマシン(ここではこれを線形サポートベクターマシンと呼ぶことにしよう)にカーネル法を適用することで柔軟性を格段に向上させたものであり, 現在でも現役で使用されている, カーネル法の応用例のひとつでもある. しかし, カーネルサポートベクターマシンには, 学習や推論にかかる時間が学習データ数に多項式的に比例してしまうという決定的な弱点がある. 具体的には, 学習データの数を N とすると, 学習に要するステップ数は $O(N^3)$, 推論に要するステップ数は $O(N)$ である. 繰り返しになるが N は学習データの数である. これは, データの質よりも量の多さをよとする機械学習の分野において, 致命的とも言える弱点である. この恐ろしさがお分かり頂けるであろうか?

想像して頂きたい. 何らかのクラス分類器を作成したところ, 性能があまり出なかったとしよう. そうなったとき, 学習データセットの拡充は間違いなく優先度の高い検討事項である. しかし, 先に述べた通り, カーネルサポートベクターマシンの場合は学習データの数を増やせば増やすほど学習や推論にかかる時間が増えてしまう. 学習時間が増えるのはまだ許容範囲内であろうが, 推論の時間まで増加してしまうのは許容し難いのではないだろうか. もちろんケースバイケースではあるのだが, 実応用では推論に費やせる計算量の上限は事前に決められている場合が多く, 推論の計算量の増加は致命的となるケースが多い. つまり, カーネルサポートベクターマシンの場合は安易に学習データセットの追加ができないのである.

ちなみに, この弱点は線形サポートベクターマシンでは起こりえない. 線形のサポートベクターマシンの場合, 推論時間は特徴量数(入力データの次元数)のみに依存し, 学習データセットの数には依存しない. 言い換えれば, この弱点はカーネル法によって高い柔軟性を得たことに由来する性質なのである. とはいえ, 線形サポートベクターマシンでは性能が出せないケースが非常に多いのも事実であるから, 何とかしてカーネルサポートベクターマシンを高速化し, 大量の学習データセットを

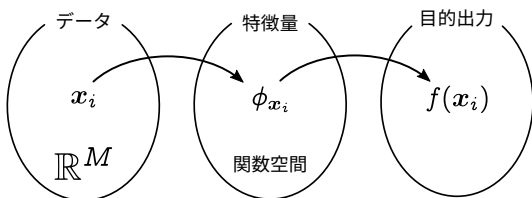


図1 カーネル法と特徴抽出

適用したいものである。

この弱点を解消するために考え出されたのが、カーネル関数の有限次元近似であり、その中で私の知る限り最も成功しているのが RFF である。

これをキチンと説明するために、まずはカーネル法について簡単にご説明していこう。

2 カーネルとは？

本節ではカーネル法の概要を、数学的証明を用いずに説明してみたいと思う。本節でカーネル法を説明する主な目的は、次節で行う RFF の解説において事前に記号の理解を持って頂くことである。厳密さが多少なりとも犠牲になるとについてはご容赦頂きたい。

やや深層学習に寄ったものの見方ではあるが、機械学習は主に「データから特徴量を抽出する」ステップと「その特徴量を用いてタスクを行う」ステップに大別される。図1を参照されたい。以下ではこの2つのステップを踏むことでカーネル法を概観してみよう。

2.1 無限次元の特徴抽出と特徴関数

入力データの集合を $\mathcal{D} \triangleq \{\mathbf{x}_i \in \mathbb{R}^M \mid i = 1, 2, \dots, M\}$ とおく。カーネル法は、入力データ \mathbf{x} から特徴量として関数、すなわち無限次元の特徴量を抽出する手法である。このデータから無限次元の特徴量を抽出する関数、これはそれを特徴関数と呼ぶこととするが、これを $\phi(\mathbf{x}, \mathbf{p})$ で表すこととしよう。具体例を挙げるとすれば、例えば

$$\phi(\mathbf{x}, \mathbf{p}) \triangleq \exp\left(-\frac{\|\mathbf{x} - \mathbf{p}\|^2}{2\sigma^2}\right), \quad (1)$$

である。これは有名なカーネル関数である RFB カーネル (radial basis function kernel) に対応する特徴関数である。

NOTE: 後に説明するが、カーネル関数と特徴関数は全くの別物であるので注意されたい。カーネル関数は特徴関数同士の内積として定義されるもので、特徴関数とは立場上全くの別物である。ただ、RBF カーネルを含む有名どころのカーネル関数は結果として再生性をみだすゆえに特徴関数とカーネル関数の関数形が一致し、説明する側としてはややこしいのだが...^^;

注意して欲しいのは、特徴関数は入力データ \mathbf{x} を受け取り、 \mathbf{p} に関する関数を返す関数である。そういう意味で、特徴関数は $\phi(\mathbf{x}, \mathbf{p})$ と書くよりも、 $\phi_{\mathbf{x}}(\mathbf{p})$ と書いた方が分かり良いであろう。やや通常の関数表記から逸脱するが、以下では $\phi_{\mathbf{x}}(\mathbf{p})$ の表記を用いることとする。また、関数であることを特に強調したい場合は、引数の \mathbf{p} を省略して $\phi_{\mathbf{x}}$ と書く。

NOTE: 余談ではあるが、関数である $\phi_{\mathbf{x}}$ を特徴「量」と呼称して良いものであろうか？これはやや疑義を生ずる点であろうが、著者は量と呼称して問題ないと考えている。関数は L^2 空間に代表されるベクトル空間を構成する要素であるから、間違いなくベクトル「量」である。

2.2 特徴関数からのタスク実施

さて、データ \mathbf{x} から特徴量 $\phi_{\mathbf{x}}$ が得られたところで、次はこの特徴量を用いて回帰や分類などのタスクを行うことを考えていこう。タスクの実施の仕方はタスクによって様々なのだが、最も一般的なアプローチは

$$f(\mathbf{x}) = \int \cdots \int_{\mathbb{R}^M} w(\mathbf{p}) \phi_{\mathbf{x}}(\mathbf{p}) d\mathbf{p}_1 \cdots d\mathbf{p}_M, \quad (2)$$

という形で関数を構築する方法である。回帰がしたいのであれば $f(\mathbf{x})$ が目的の曲線を追従するように、分類がしたいのであれば $f(\mathbf{x})$ が対数尤度を出力するようにすれば良い。ただし $w(\mathbf{p})$ は学習パラメータに相当する関数である。この関数を適切に学習させることで目的のタスクを実現するのがカーネル法である。

NOTE: 式(2)が分かりづれえ！と思ったその貴方。是非とも関数解析を習得されることをお薦めする。式(2)はただの関数同士の「内積」であり、 $f(\mathbf{x}) = \langle w, \phi_{\mathbf{x}} \rangle$ と表記することもできる。つまりやっていることは線形回帰と何も変わらず、ただ扱う空間がベクトル空間から関数空間に変わっただけのことである。

もちろん、最適化パラメータが関数であるような最適化問題は数値計算問題として解きにくいので、様々なテクニックを駆使してこれを有限次元の最適化問題に落とし込み、その最適化を問題を解くことで学習を行うのである。

メデタシ、メデタシ...

2.3 いやいや、カーネル関数はどこよ？

え、カーネル関数はいつ出てくるのかって？実はカーネル関数は、最適化問題をキチンと定式化し、その式をコネクリ回さないで出てこない。本節ではその部分の理解にあまり体力を使いたくないので、先人が数式をコネクリ回した結果、以下が得られたということだけ指摘しておくに留めておく。詳細については赤穂 [3] や Bishop [4] などを参照されたい。

1. 関数 w の最適化は有限次元のベクトルの最適化に落とし込むことが出来る。ただしそのベクトルの次元は学習データ数 N である。この関数の最適化問題を有限次元の最適化問題に落とし込むテクニックは、カーネルトリックと呼ばれる。
2. 上述1の最適化を解く上でも、最適化の結果を用いて推論を行う上でも、実は特徴関数 $\phi_{\mathbf{x}}$ は陽には出てこない。出てくるのは次の式

$$k(\mathbf{x}, \mathbf{y}) \triangleq \int \cdots \int \phi_{\mathbf{x}}(\mathbf{p}) \phi_{\mathbf{y}}(\mathbf{p}) d\mathbf{p}_1 \cdots d\mathbf{p}_M, \quad (3)$$

で与えられる関数 k だけである。この関数 k をカーネル関数と呼ぶ。

さて、上記には2つの重要な事実が含まれている。ひとつは、カーネルトリックによって関数次元の最適化問題を有限次元の

最適化問題に落とし込めたとはいえ、いまだに学習データ数 N と同じだけの次元を持っているという点である。これが本文書の冒頭で述べたカーネル法の弱点に直結しているのである。

もうひとつは、学習や推論を行う上で特徴関数は陽には出でず、カーネル関数しか出てこないという点である。ゆえに、一般的にはカーネル関数だけを設計し、特徴ベクトルには言及しない場合が多い。ちなみに、式 (1) にて特徴関数の具体例をお見せしたが、これに対応するカーネル関数は

$$k(\mathbf{x}, \mathbf{y}) = \exp(-\gamma \|\mathbf{x} - \mathbf{y}\|^2), \quad \gamma = -\frac{1}{4\sigma^2}, \quad (4)$$

である。これは式 (3) を直接計算することで得られる。

さて、様々な式展開や証明を省略させて頂いたが、これでカーネル法のおさらいは終わりにして、早速 RFF の定式化に移ろう。

NOTE: カーネル関数の説明が適当すぎないかって？ご指摘の通りであるが、本文書としてはカーネル法の説明は前座にすぎないので、本節でパルプを3枚も4枚も費すことはしたくない。もちろんカーネル関数がいきなり天から降ってくるのは本文書の想定読者を考えると論外である。もっとスマートな導入の仕方があれば是非ご連絡頂きたい。

3 カーネルの有限近似

3.1 実確率変数の期待値

突然ではあるが、確率変数の期待値について簡単に説明をさせて頂きたい。まずは簡単のため、集合 $\{z_1, z_2, \dots, z_I\}$ 上のいずれかの値を取る離散確率変数 Z を考える。この確率変数 Z の期待値は

$$\mathbb{E}[X] = \sum_{i=1}^I p_i z_i, \quad (5)$$

で与えられる。ただし p_i は Z が値 z_i を取る確率である。連続確率変数でも同じことが言える。確率変数 X をとある確率空間上で可測な実数値確率変数とし、 X が実数値 x を取る確率密度関数を $p(x)$ とおく。このとき確率変数 X の期待値は

$$\mathbb{E}[X] = \int_{\mathbb{R}} p(x) x \, dx, \quad (6)$$

で与えることができる。まさに離散確率変数の期待値を連続化したような式になっている。ちなみに確率変数 X が M 次元のベクトル確率変数であった場合は

$$\mathbb{E}[X] = \int \cdots \int_{\mathbb{R}} p(\mathbf{x}) \mathbf{x} \, dx_1 \cdots dx_M, \quad (7)$$

となる。

さらに、式 (6) の右辺を解析的に解くことが困難であったとしよう。このとき我々には何が出来るであろうか。そう、期待値なのであるから、実際に確率変数 X をいくつか生成してみ、その平均値をもって期待値の近似値とすることができる。これはいわゆる Monte Carlo 法であり、数式で表現すると

$$\mathbb{E}[X] \approx \frac{1}{L} \sum_{\ell=1}^L X_{\ell}, \quad (8)$$

となる。ただし X_{ℓ} は確率変数 X のサンプリングである。

何でいきなりこんな話をしたかって？マァマァ、黙って式 (6), (7) を目に焼き付けておいて頂きたい。

3.2 カーネルの Fourier 変換

さて、本題の RFF に話を戻そう。以下ではカーネル関数 $k(\mathbf{x}, \mathbf{y})$ を $k(\mathbf{x} - \mathbf{y})$ と表現できるものだけに限定しよう。多項式カーネルや Fisher カーネルなど多くの有名なカーネル関数が議論の対象外になってしまうものの、式 (4) で与えられる RBF カーネルは幸いにも $k(\mathbf{x} - \mathbf{y})$ はそのように書き表すことができる。

さて、カーネル関数 $k(\mathbf{x} - \mathbf{y})$ を Fourier 変換し、その結果を $p(\boldsymbol{\omega})$ と置いてみよう。すなわち

$$\begin{aligned} p(\boldsymbol{\omega}) &= \mathcal{F}[k](\boldsymbol{\omega}) \\ &= \int \cdots \int_{\mathbb{R}^M} k(\mathbf{z}) e^{-i\boldsymbol{\omega}^T \mathbf{z}} \, dz_1 \cdots dz_M, \end{aligned} \quad (9)$$

である。これを逆 Fourier 変換を用いて $k(\mathbf{x} - \mathbf{y}) = \cdots$ の形に書き直すと

$$\begin{aligned} k(\mathbf{x} - \mathbf{y}) &= \mathcal{F}^{-1}[p](\mathbf{x} - \mathbf{y}) \\ &= \int \cdots \int_{\mathbb{R}^M} p(\boldsymbol{\omega}) e^{i\boldsymbol{\omega}^T (\mathbf{x} - \mathbf{y})} \, d\omega_1 \cdots d\omega_M, \end{aligned} \quad (10)$$

が成り立つ。これは逆 Fourier 変換の定理から明らかであろう。

さて、式 (10) を見て思うことはないだろうか。そう、期待値である。つまり p を確率密度関数とみなせば

$$k(\mathbf{x} - \mathbf{y}) = \mathbb{E}_{\boldsymbol{\omega} \sim p} \left[e^{i\boldsymbol{\omega}^T (\mathbf{x} - \mathbf{y})} \right], \quad (11)$$

と表すことができる。ただし、上式の期待値は $\boldsymbol{\omega}$ を確率変数とみなし、さらにその確率分布は確率密度関数 $p(\boldsymbol{\omega})$ に従うとみなしていることに注意せよ。期待値記号 \mathbb{E} の添字に $\boldsymbol{\omega} \sim p$ を書いたのは、そのことを明記するためである。

先に指摘した通り、我々は期待値計算を Monte Carlo 法で近似できることを知っている。確率密度関数 $p(\boldsymbol{\omega})$ にしたがって確率変数 $\boldsymbol{\omega}$ をサンプリングした値を $\{\boldsymbol{\omega}_{\ell}\}_{\ell=1}^L$ とすると、

$$\begin{aligned} k(\mathbf{x} - \mathbf{y}) &= \mathbb{E}_{\boldsymbol{\omega} \sim p} \left[e^{i\boldsymbol{\omega}^T (\mathbf{x} - \mathbf{y})} \right] \\ &\approx \frac{1}{L} \sum_{\ell=1}^L e^{i\boldsymbol{\omega}_{\ell}^T (\mathbf{x} - \mathbf{y})}, \end{aligned} \quad (12)$$

が成り立つのである。これにより、カーネル関数を、Fourier 変換と乱数を用いて近似できたことになる。これが RFF の本質である。

さて、興奮に任せて大切なポイントをひとつ取り落としてしまったことは読者もお気づきであろう。式 (11) が成立するためには、関数 $p(\boldsymbol{\omega})$ が確率密度関数と同格でなければならない。すなわち、可積分かつ非負の値を取る関数でなければならない。しかし驚くべきことに、カーネル関数が正定値関数である限り、その Fourier 変換は可積分かつ非負であることが証明できる。これは非常に大事な事実であるから、定理としてまとめておく。ただし、議論の流れを妨げないよう証明は本文書の最後の節で行う。

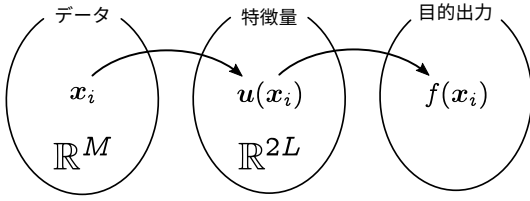


図2 RFF と特徴抽出

定理 3.1 (平行移動不変な正定値カーネルの Fourier 変換)

正定値カーネル関数 $k(\mathbf{x}, \mathbf{y})$ が平行移動不変，すなわち $k(\mathbf{x}, \mathbf{y}) = k_\Delta(\mathbf{x} - \mathbf{y})$ をみたす関数 k_Δ が存在したとする．このとき，カーネル関数 $k_\Delta(\mathbf{x} - \mathbf{y})$ の Fourier 変換は L_1 ノルムの意味で可積分であり，かつ非負関数である．

さて，最後に近似したカーネル関数をどのように使うかについて説明しよう．一般に，カーネル関数は正定値関数であるから，式 (12) の最右辺の虚部は不要である．したがって

$$k(\mathbf{x} - \mathbf{y}) \approx \frac{1}{L} \sum_{\ell=1}^L e^{i\omega_\ell^\top (\mathbf{x} - \mathbf{y})} \quad (13)$$

$$= \frac{1}{L} \sum_{\ell=1}^L \cos(\omega_\ell^\top (\mathbf{x} - \mathbf{y})) \quad (14)$$

$$= \frac{1}{L} \sum_{\ell=1}^L [\cos(\omega_\ell^\top \mathbf{x}) \cos(\omega_\ell^\top \mathbf{y}) + \sin(\omega_\ell^\top \mathbf{x}) \sin(\omega_\ell^\top \mathbf{y})], \quad (15)$$

を得る．ただし式 (13) から式 (14) へは Euler の公式を，式 (14) から式 (15) へはかの懐しき三角関数の加法定理を用いた．さてここでベクトル

$$\mathbf{u}(\mathbf{x}) \triangleq \frac{1}{\sqrt{L}} \begin{pmatrix} \cos(\omega_1^\top \mathbf{x}) \\ \vdots \\ \cos(\omega_L^\top \mathbf{x}) \\ \sin(\omega_1^\top \mathbf{x}) \\ \vdots \\ \sin(\omega_L^\top \mathbf{x}) \end{pmatrix} = \begin{pmatrix} \cos(\mathbf{W}\mathbf{x}) \\ \sin(\mathbf{W}\mathbf{x}) \end{pmatrix} \quad (16)$$

を導入しよう．ただし行列 $\mathbf{W} = (\omega_1, \dots, \omega_L)^\top$ であり，上式の最右辺ではベクトルの余弦関数を $\cos \mathbf{x} = (\cos x_1, \dots, \cos x_L)^\top$ と定義して使用している． \sin についても同様である．すると式 (15) の右辺は

$$\mathbf{u}(\mathbf{x})^\top \mathbf{u}(\mathbf{y}), \quad (17)$$

と非常にシンプルに書き改めることができる．すなわち

$$k(\mathbf{x} - \mathbf{y}) \approx \mathbf{u}(\mathbf{x})^\top \mathbf{u}(\mathbf{y}), \quad (18)$$

と近似することができたのである．

そもそもカーネル関数は，特徴量の内積であったことを思い出してほしい．式 (18) は，カーネル関数を近似することと，特徴量として $\phi_{\mathbf{x}}$ ではなく $\mathbf{u}(\mathbf{x})$ を使うことが等価だと言っているのである．つまり図 2 のような構図が成り立つ．

今，特徴量はもはや関数ではなく，ただの $2L$ 次元の実ベクトルである．したがって，回帰や分類などのタスクを実行するのにミョーチクリンなトリックはもはや必要ない．有限次元の線形回帰なり線形サポートベクターマシンなりを適用すれば良いのである．線形サポートベクターマシンの学習および推論は，カーネルサポートベクターマシンと比較すると桁違いに速い．しかも特徴量は，高々有限次元の近似とはいえ，カーネル関数の柔軟性を残している．このようにしてカーネル法はその柔軟性を残しつつ大規模データに適用しうるだけの高速性を手にすることができるのである．

3.3 RBF カーネルの場合

さて，実際に RFF を運用するためには，カーネル関数を何らかの関数に定め，そのカーネル関数に対して確率密度関数 $p(\omega)$ を算出しなければならない．RFF はカーネル関数が $k(\mathbf{x} - \mathbf{y})$ と表現できるものに限られるから，ここでは RBF カーネルの場合について考察しよう．

とはいえ，確率密度関数 $p(\omega)$ の導出は式 (9) にしたがって Fourier 変換を計算するだけである．直接計算だけであるので証明は省略するが，

$$p(\omega) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{\|\omega\|^2}{2\sigma^2}\right), \quad (19)$$

となる．つまり確率変数 ω は正規分布するのである．したがって確率変数 ω は正規分布からサンプリングすれば良いということになる．

以上をまとめると，RBF ベースの RFF は次の手順で実施すれば良い．

1. RFF の自由度 $L \in \mathbb{Z}^+$ を適当な数に設定する．これはカーネル関数を近似する際のサンプリング数である．
2. 行列の各要素が正規分布にしたがうランダム行列 $\mathbf{W} \in \mathbb{R}^{L \times M}$ を生成する．
3. 各学習データ \mathbf{x}_i に式 (16) の変換を施し，それを改めて学習データセットとみなす．
4. 変換後のデータセットに対し，線形回帰や線形サポートベクターマシンなどのアルゴリズム f を適用する．
5. 推論時には，推論データ \mathbf{x} に対してもデータ変換式 (16) を適用して $\mathbf{u}(\mathbf{x})$ を算出し，その後に推論 $f(\mathbf{u}(\mathbf{x}))$ を行う．

3.4 数値実験：RBF カーネルの RFF 近似

本節では，RFF を用いて RBF カーネルが近似される様子を実際に確認してみよう．まずは近似前の RBF カーネルを

$$k(x_1, x_2) \triangleq \exp\left(-\frac{\|x_1 - x_2\|^2}{2}\right), \quad (20)$$

とおくとしよう．このとき対応する RFF の確率密度関数は

$$p(\omega) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\|\omega\|^2}{2}\right), \quad (21)$$

となるので，RFF 近似した RBF カーネルは

$$k_{\text{rff}}(x_1, x_2) \triangleq \mathbf{u}(x_1)^\top \mathbf{u}(x_2), \quad \omega_1, \dots, \omega_L \sim p(\omega), \quad (22)$$

となる。これを試みに $L = 10, 100, 1000$ として描画してみたのが図 3 である。RFF の次元を上げるにつれて RBF カーネルに近づいていくことが分かるであろう。

3.5 深層学習との関係

最後に、RFF を用いたサポートベクターマシンと深層学習との関係性について述べておく。RFF は入力データ \mathbf{x} の特徴量を式 (16) で定める手法であった。これを線形サポートベクターマシンと組み合わせてみよう。

問題設定としては、教師付き学習データを用いて C クラスのクラス分類器を作成することを目標としよう。学習データセットは $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ とする。ただし $\mathbf{x}_n \in \mathbb{R}^M$, $y_n \in \{1, 2, \dots, C\}$ であり、 \mathbf{x}_n が学習データに、 y_n が教師データに対応する。

さて、線形サポートベクターマシンの判別式は線形関数であるから、入力データ \mathbf{x} に対する各クラスの対数尤度は

$$f_c(\mathbf{x}) = \mathbf{a}_c^T \mathbf{x} + b_c, \quad (23)$$

で表現される。ベクトル \mathbf{a}_c およびスカラー b_c は学習パラメータであり、各クラスの対数尤度を上手く表現できるように学習される。

さて、少しだけ表現を変更しよう。まず、すべてのクラスをまとめて記述するために、 $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_C)^T$, $\mathbf{b} = (b_1, \dots, b_C)$ とおく。さらに、RFF の適用を前提とすれば、入力データは \mathbf{x} ではなく $\mathbf{u}(\mathbf{x})$ に置き換えなければならない。すると式 (23) は

$$\mathbf{f}(\mathbf{x}) = \mathbf{A} \begin{pmatrix} \cos(\mathbf{W}\mathbf{x}) \\ \sin(\mathbf{W}\mathbf{x}) \end{pmatrix} + \mathbf{b}, \quad (24)$$

となる。

これを深層学習の立場から見ると、ちょっと特殊な活性化関数を用いた、2 層の全結合層を持つニューラルネットワークとみなすことができる。図 4 を参照されたい。

もちろん、学習の仕方はサポートベクターマシンと深層学習で全く異なることに注意されたい。サポートベクターマシンは区分的線形な目的関数を制約付き 2 次計画問題に帰着させているのに対し、深層学習ではクロスエントロピーを確率的勾配法で最小化させている。

4 数値実験

GitHub に RFF の実装およびそれを用いたいいくつかの計算例を掲載してある。以下のレポジトリを参照されたい。

<https://github.com/tiskw/>

Random-Fourier-Features/tree/master/
examples/rff_svc_for_mnistsubsection

5 定理の証明

本節では、本文書で証明せずに残してあった定理を証明する。前節までは、原論文に合わせるために角周波数ベースの Fourier 変換を用いて議論をしていたが、以下では数学的シンプルさを優先して周波数ベースの Fourier 変換を用いることとする。

Fourier 変換と Parseval の等式

まずは Fourier 変換について簡単におさらいしておこう。本節では Fourier 変換の定義として

$$\mathcal{F}[f](\xi) \triangleq \int_{-\infty}^{\infty} f(x) e^{-2\pi i x \xi} dx, \quad (25)$$

を採用する。この時、逆 Fourier 変換は

$$\mathcal{F}^{-1}[\hat{f}](x) \triangleq \int_{-\infty}^{\infty} \hat{f}(\xi) e^{2\pi i x \xi} d\xi, \quad (26)$$

となる。原論文で用いられている角周波数ベースの Fourier 変換とはスカラー倍の違いがあるだけで、本質的には同一である。多変数関数 $f(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^M$ の Fourier 変換および逆 Fourier 変換は

$$\mathcal{F}[f](\xi) \triangleq \int_{-\infty}^{\infty} f(\mathbf{x}) e^{-2\pi i \mathbf{x}^T \xi} d\mathbf{x}, \quad (27)$$

$$\mathcal{F}^{-1}[\hat{f}](\mathbf{x}) \triangleq \int_{-\infty}^{\infty} \hat{f}(\xi) e^{2\pi i \mathbf{x}^T \xi} d\xi, \quad (28)$$

となる。ただし記号 $d\mathbf{x}$ は $dx_1 dx_2 \dots dx_M$ の略記であることに注意せよ。

NOTE: Lebesgue 測度論の慣習に従えば、この $d\mathbf{x}$ は本来 $\mu(d\mathbf{x})$ あるいは $d\mu(\mathbf{x})$ と表記されるべきものである。ただし μ は Lebesgue 測度である。だが表記の冗長さと、測度論に立ち入り結論に辿り着くまでの時間が長くなることで読者の興味が薄れることを恐れ、ここでは上記の略記を用いた。ゴマカシだという誇りは受けて然るべきである。より良いアプローチがあればぜひご連絡頂きたい。

オイオイ、式 (25) が定義できる関数 f なんてかなり少ないかと思ったその読者、非常に良い眼をお持ちである。式 (25) の収束性については本文書末尾の付録で簡単に言及するが、応用上必要な関数に対しては、様々なテクニックを駆使することで収束させることができるので安心して頂きたい。

次に Parseval の等式 (*Parseval's identity*) について見ていこう。Parseval の等式と言うと公式 $\|f\|^2 = \|\hat{f}\|^2$ を指すことが多いが、本文書ではより一般的な形の定理を紹介する。

定理 5.1 (Parseval の等式)

任意の Fourier 変換可能な関数 f, g およびその Fourier 変換 \hat{f}, \hat{g} に対して $\langle f, g \rangle = \langle \hat{f}, \hat{g} \rangle$ が成り立つ。

ここでは Parseval の等式の証明は省略する。証明が気になる読者は（そうあるべきである！）本文書末尾の付録を参照されたい。聡明な読者には無用の言であろうが、 $g = f$ とすれば $\|f\|^2 = \|\hat{f}\|^2$ を得る。

定理の証明

定理を再掲し証明しよう。

定理 5.2 (平行移動不変な正定値カーネルの Fourier 変換)

正定値カーネル関数 $k(\mathbf{x}, \mathbf{y})$ が平行移動不変、すなわち $k(\mathbf{x}, \mathbf{y}) = k_{\Delta}(\mathbf{x} - \mathbf{y})$ をみたす関数 k_{Δ} が存在したとする。このとき、カーネル関数 $k_{\Delta}(\mathbf{x} - \mathbf{y})$ の Fourier 変換は L_1 ノルムの意味で可積分であり、かつ非負関数である。

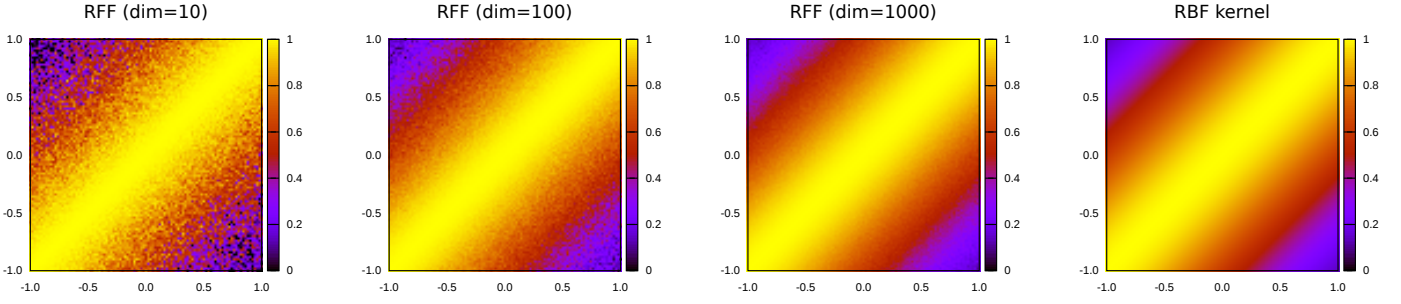


図3 RBF カーネルの RFF 近似

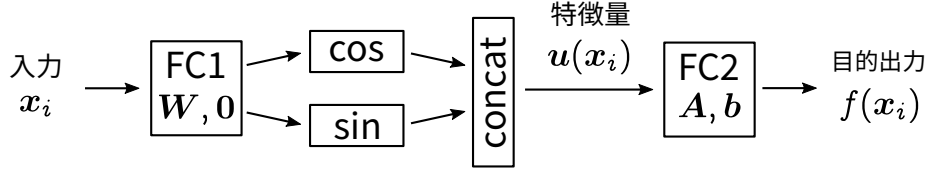


図4 RFF とニューラルネットワーク

証明： 以下、 \mathbf{x}, \mathbf{y} の次元を M とし、 $\mathcal{X} = \mathbb{R}^M$ とおく．すなわち $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ である．カーネル関数 k の対称性より

$$k_{\Delta}(\mathbf{x}) = k(\mathbf{x}/2, -\mathbf{x}/2) = k(-\mathbf{x}/2, \mathbf{x}/2) = k_{\Delta}(-\mathbf{x}) \quad (29)$$

が成り立つ．これを用いれば k_{Δ} の Fourier 変換は実数値関数であることがわかる．なぜならば k_{Δ} の Fourier 変換の複素共役を計算してみると

$$\mathcal{F}[k_{\Delta}]^*(\boldsymbol{\xi}) = \int_{\mathcal{X}} k_{\Delta}(\mathbf{x}) e^{2\pi i \mathbf{x}^T \boldsymbol{\xi}} d\mathbf{x} \quad (30)$$

ここで $\mathbf{x} = -\mathbf{z}$ と置換すれば

$$= \int_{\mathcal{X}} k_{\Delta}(-\mathbf{z}) e^{-2\pi i \mathbf{z}^T \boldsymbol{\xi}} d\mathbf{z} \quad (31)$$

$$= \int_{\mathcal{X}} k_{\Delta}(\mathbf{z}) e^{-2\pi i \mathbf{z}^T \boldsymbol{\xi}} d\mathbf{z} = \mathcal{F}[k_{\Delta}](\boldsymbol{\xi}), \quad (32)$$

となり，複素共役をとる前の Fourier 変換と一致する．したがって $\mathcal{F}[k_{\Delta}](\boldsymbol{\xi})$ は少なくとも実数値関数である．次に，カーネル関数 k が正定値関数であることから，任意の実数値関数 y について

$$\iint_{\mathcal{X} \times \mathcal{X}} y(\mathbf{x}_1) k(\mathbf{x}_1, \mathbf{x}_2) y(\mathbf{x}_2) d\mathbf{x}_1 d\mathbf{x}_2 \geq 0, \quad (33)$$

が成り立つ．上式の左辺を変形すると

$$\iint_{\mathcal{X} \times \mathcal{X}} y(\mathbf{x}_1) k(\mathbf{x}_1, \mathbf{x}_2) y(\mathbf{x}_2) d\mathbf{x}_1 d\mathbf{x}_2 \quad (34)$$

$$= \int_{\mathcal{X}} y(\mathbf{x}_1) \langle k(\mathbf{x}_1, \cdot), y \rangle d\mathbf{x}_1 \quad (35)$$

ここで内積 $\langle k(\mathbf{x}_1, \cdot), y \rangle$ に対してパーセバルの不等式を用いる．新たに $\hat{k}_{\mathbf{x}_1}(\boldsymbol{\xi}_2) = \mathcal{F}[k(\mathbf{x}_1, \cdot)](\boldsymbol{\xi}_2)$ とおくと

$$= \int_{\mathcal{X}} y(\mathbf{x}_1) \langle \hat{k}_{\mathbf{x}_1}, \hat{y} \rangle d\mathbf{x}_1 \quad (36)$$

$$= \iint_{\mathcal{X} \times \mathcal{X}} y(\mathbf{x}_1) \hat{k}_{\mathbf{x}_1}(\boldsymbol{\xi}_2) \hat{y}^*(\boldsymbol{\xi}_2) d\mathbf{x}_1 d\boldsymbol{\xi}_2 \quad (37)$$

ここで見易さのために関数 $\hat{k}_{\mathbf{x}_1}(\boldsymbol{\xi}_2) = l_{\boldsymbol{\xi}_2}(\mathbf{x}_1)$ とおくと

$$= \iint_{\mathcal{X} \times \mathcal{X}} y(\mathbf{x}_1) l_{\boldsymbol{\xi}_2}(\mathbf{x}_1) \hat{y}^*(\boldsymbol{\xi}_2) d\mathbf{x}_1 d\boldsymbol{\xi}_2 \quad (38)$$

$$= \int_{\mathcal{X}} \langle y, l_{\boldsymbol{\xi}_2} \rangle \hat{y}^*(\boldsymbol{\xi}_2) d\boldsymbol{\xi}_2 \quad (39)$$

ここで再び $\langle y, l_{\boldsymbol{\xi}_2} \rangle$ に対してパーセバルの不等式を用いる．新たに $\hat{l}_{\boldsymbol{\xi}_2}(\boldsymbol{\xi}_1) = \mathcal{F}[l_{\boldsymbol{\xi}_2}](\boldsymbol{\xi}_1)$ とおくと

$$= \int_{\mathcal{X}} \langle \hat{y}, \hat{l}_{\boldsymbol{\xi}_2} \rangle \hat{y}^*(\boldsymbol{\xi}_2) d\boldsymbol{\xi}_2 \quad (40)$$

$$= \iint_{\mathcal{X} \times \mathcal{X}} \hat{y}(\boldsymbol{\xi}_1) \hat{l}_{\boldsymbol{\xi}_2}(\boldsymbol{\xi}_1) \hat{y}^*(\boldsymbol{\xi}_2) d\boldsymbol{\xi}_1 d\boldsymbol{\xi}_2, \quad (41)$$

となる．ここで $\hat{l}_{\boldsymbol{\xi}_2}(\boldsymbol{\xi}_1)$ を真面目に計算してみると

$$\hat{l}_{\boldsymbol{\xi}_2}(\boldsymbol{\xi}_1) = \iint_{\mathcal{X} \times \mathcal{X}} k(\mathbf{x}_1, \mathbf{x}_2) e^{-2\pi i (\mathbf{x}_1^T \boldsymbol{\xi}_1 - \mathbf{x}_2^T \boldsymbol{\xi}_2)} d\mathbf{x}_1 d\mathbf{x}_2 \quad (42)$$

$$= \iint_{\mathcal{X} \times \mathcal{X}} k_{\Delta}(\mathbf{x}_1 - \mathbf{x}_2) e^{-2\pi i (\mathbf{x}_1^T \boldsymbol{\xi}_1 - \mathbf{x}_2^T \boldsymbol{\xi}_2)} d\mathbf{x}_1 d\mathbf{x}_2 \\ = \delta(\boldsymbol{\xi}_1 - \boldsymbol{\xi}_2) \hat{k}(\boldsymbol{\xi}_1), \quad (43)$$

となる．これを式 (41) に戻すと

$$= \iint_{\mathcal{X} \times \mathcal{X}} \hat{y}(\boldsymbol{\xi}_1) \delta(\boldsymbol{\xi}_1 - \boldsymbol{\xi}_2) \hat{k}(\boldsymbol{\xi}_1) \hat{y}^*(\boldsymbol{\xi}_2) d\boldsymbol{\xi}_1 d\boldsymbol{\xi}_2 \quad (44)$$

$$= \int_{\mathcal{X}} \hat{y}(\boldsymbol{\xi}) \hat{k}(\boldsymbol{\xi}) \hat{y}^*(\boldsymbol{\xi}) d\boldsymbol{\xi} \quad (45)$$

$$= \int_{\mathcal{X}} \|\hat{y}(\boldsymbol{\xi})\|^2 \hat{k}(\boldsymbol{\xi}) d\boldsymbol{\xi} \geq 0, \quad (46)$$

となる．これが任意の関数 y について成立するから， $\hat{y}(\boldsymbol{\xi})$ が定数関数の場合も成立しなければならない．したがって

$$\int_{\mathcal{X}} \hat{k}(\boldsymbol{\xi}) d\boldsymbol{\xi} \geq 0, \quad (47)$$

が成り立つ必要がある．したがって \hat{k} は L_1 ノルムの意味で可積分である．同様に関数 $\hat{y}(\boldsymbol{\xi})$ が限りなくデルタ関数に近い形

状の関数に対しても式 (41) が成立しなければならないから、任意の ξ に対して $\widehat{k}(\xi) \geq 0$ が成り立つ必要がある。したがって $\widehat{k}(\xi) \geq 0$ は非負関数である。 ■

付録：Fourier 変換

本節では、定理証明の際に省略した Fourier 変換および Parseval の等式の解説を行う。本節でも数学的シンプルさを優先して周波数ベースの Fourier 変換を用いることとする。

Fourier 変換に関する補足

Fourier 変換の定義としては本文の式 (25) で十分だが、より良い表記、より深い理解のために、Fourier 変換を関数空間の内積を用いて表現してみよう。まず関数 $\eta_\xi(x)$ を

$$\eta_\xi(x) \triangleq e^{2\pi i x \xi}, \quad (48)$$

と定義する。次に関数の内積 $\langle \cdot, \cdot \rangle$ を最も標準的な内積演算

$$\langle f, g \rangle \triangleq \int_{-\infty}^{\infty} f(x) g^*(x) dx, \quad (49)$$

で定めるとしよう。ただし $g^*(x)$ は関数 $g(x)$ の複素共役である。このとき式 (25) で表現でされる Fourier 変換は

$$\mathcal{F}[f](\xi) = \langle f, \eta_\xi \rangle, \quad (50)$$

と表すことができる。このとき関数群 $\{\eta_\xi(x) | \xi \in \mathbb{R}\}$ は正規直交基底と非常に似た性質を持つ。すなわち

$$\langle \eta_{\xi_1}, \eta_{\xi_2} \rangle = \delta(\xi_1 - \xi_2), \quad (51)$$

が成り立つ。ただし δ は Dirac のデルタ関数 (Dirac delta function) である。つまり式 (51) は $\xi_1 = \xi_2$ のときのみ値を持つ。これは通常のベクトル空間の基底 $\{e_1, e_2, \dots\}$ が持つ性質

$$e_i^T e_j = \begin{cases} 1, & i = j, \\ 0, & i \neq j, \end{cases} \quad (52)$$

と非常に良く似ている。したがって Fourier 変換は正規直交基底との内積、つまるところただの基底交換のようなもの、と見ることもできるのである。

NOTE: Fourier 変換の定義を初めて見た若かりし頃、諸君らも「何で Fourier 変換が逆回転 $e^{-2\pi i x \xi}$ で逆 Fourier 変換が順回転 $e^{2\pi i x \xi}$ なんだよ、フツーは逆だろ！」とキレかかった経験があるのではないだろうか。上述の内積による Fourier 変換の理解はこの疑問を解消してくれる。そう、基底 $\eta_\xi(x)$ はちゃんと順回転 $e^{2\pi i x \xi}$ で定義されていたのだ。しかし内積を計算する際に共役をとる必要があり、そこで回転方向が反転したのである。少しはスッキリしたであろうか。

Fourier 変換の収束性について

Parseval の等式の証明に移る前に、一部の読者が気になっているであろう Fourier 変換 (25) の収束性について簡単に述べておこう。つまり、式 (25) が定義できるのはどんな関数 f か？という問題である。結論から言えば、式 (25) が収束する条件は

$$\int_{-\infty}^{\infty} |f(x)| dx < \infty, \quad (53)$$

である。ハッキリ言って、上式を満たす関数は応用上の観点から見ても非常に少ない。例えば三角関数や指数関数ですら上式の制約は満たさない。もちろん世の中には病的な関数が沢山あるゆえ、真に全ての関数に対して Fourier 変換を定義せよというのは困難であるのだが、とはいえ、少なくとも応用上必要な関数に対しては Fourier 変換ができないと困ってしまう。この問題に対するアプローチは様々あるが、著者の知る限りでは、

- (1) 扱う関数はすべて有限の台 (support) を持つと仮定する方法。すなわち、関数が値を持つ範囲を十分に広い有限な範囲に限定してしまうという方法。関数の台が有限であれば、応用上必要な関数はほぼすべて積分可能となる。
- (2) 両側 Laplace 変換の特殊な場合 ($s = i\omega$) として Fourier 変換を定義する方法である。実は両側 Laplace 変換は、解析接続を用いた巧みな仕組みにより、非常に幅広い関数に対して定義可能である。

の 2 つのアプローチがある。初学者は (1) の理解で良いであろう。世の中の多くの人の理解は以下の (1) の方法で理解していると思われる。ちなみに著者は理論的にも便利な (2) で理解している。興味があれば調べてみると良い。

Parseval の等式

次に Parseval の等式 (Parseval's identity) について見ていこう。一般に Parseval の等式と言うと $\|f\|^2 = \|\widehat{f}\|^2$ を指すことも多いが、本文書ではより一般的な形の定理を証明する。Fourier 変換が正規直交基底を用いた基底の交換であることが分かれば、ほぼ自明な定理である。

定理 5.3 (Parseval の等式)

任意の Fourier 変換可能な関数 f, g およびその Fourier 変換 \widehat{f}, \widehat{g} に対して $\langle f, g \rangle = \langle \widehat{f}, \widehat{g} \rangle$ が成り立つ。

証明： 証明は直接計算による。

$$\begin{aligned} \langle \widehat{f}, \widehat{g} \rangle &= \int_{-\infty}^{\infty} \langle f, \eta_\xi \rangle \langle g, \eta_\xi \rangle d\xi \\ &= \iiint_{-\infty}^{\infty} f(x_1) g(x_2) \eta_\xi(x_1) \eta_\xi(x_2) dx_1 dx_2 d\xi \end{aligned}$$

ここで $\eta_\xi(x_1) \eta_\xi(x_2) = \eta_{x_1}(\xi) \eta_{x_2}(\xi)$ に注意すれば

$$= \iint_{-\infty}^{\infty} f(x_1) g(x_2) \langle \eta_{x_1}, \eta_{x_2} \rangle dx_1 dx_2$$

関数 $\eta_x(\xi)$ の正規直交性に注意すれば $\langle \eta_{x_1}, \eta_{x_2} \rangle$ は $\delta(x_1 - x_2)$ であるから

$$= \int_{-\infty}^{\infty} f(x) g(x) dx = \langle f, g \rangle,$$

となり、題意が成り立つ。 ■

ここまで読み進めた読者には指摘するまでもないであろうが、上記定理の特殊な場合として $f = g$ とおけば、良く知られた $\|f\|^2 = \|\widehat{f}\|^2$ を得る。

本文書では、非常に簡単ではありますが、Random Fourier Features の解説を試みてみました。個人的には、簡単のためにデータを 1 次元にして説明しようかとも考えたのですが、さすがに実応用の段階でつまづくであろうと思い、最初から多次元のデータを想定して書いてみました。この点についてはもしかしたら将来的に修正するかもしれません。

Random Fourier Features の拡張としては、2017 年に Google より ORF (*Orthogonal Random Features*) および SORF (*Structured ORF*) が提案されています。これは RFF のランダム行列 \mathbf{W} が直行行列でないことに注目し、直交性を強制的に持たせることでより高い近似効率を狙ったものです。私も手元で実装し、いくつかの有名なデータセットに対して ORF および SORF を試しましたが、確かにカーネルの近似精度は高いのですが、各データセットの推論精度という意味では RFF と大きな差はありませんでした。どうやらカーネルの近似精度と推論精度は単純には相関しないようです。まだ仮説の域を出ませんが、おそらく有限次元近似が良い正則化として機能しているのでしょう。本文書で ORF や SORF に触れていないのは、それが主な理由です。

私の数学の師は常に「人生で一度も証明を与えたことのない定理を自ら使用するのには数学的犯罪である」と私に教えて下さいました。これは教育的な配慮だけでは決してないと思います。内部を理解せずに扱うことによる誤用、その招く惨めな結末。数学だけに限った話ではないでしょう。近年の機械学習の研究者やエンジニアはこれをよくよく心に留めおくべきです。もちろん私自身を含めて、です。

最後に、Toyota Research Institute Advanced Development, Inc. の石原昌弘さんには本文書の誤植を指摘して頂きました。この場を借りて厚く御礼申し上げます。また、私の数学的活動は、2017 年に逝去された恩師、山下弘一郎先生や、大学および大学院で私の指導教官を担当して下さいた早川朋久准教授をはじめ、数学で私と関わりを持ったすべての方々のおかげで成り立っています。そして、数学的活動の以前に、そもそも私の生は両親によって与えられ、妻によって支えられています。

参考文献

- [1] A. Rahimi and B. Recht, “Random Features for Large-Scale Kernel Machines”, Neural Information Processing Systems, 2007.
- [2] F. X. Yu, A. T. Suresh, K. Choromanski, D. H. Rice and S. Kumar, “Orthogonal Random Features”, Neural Information Processing Systems, 2016.
- [3] 赤穂昭太郎, 「カーネル多変量解析」, 岩波書店, 2008.
- [4] C. Bishop, “Pattern Recognition and Machine Learning”, Springer, 2010.