

1 機械学習との出会い

前職(京大化研)から現職(滋賀大学)に異動する直前の 2015 年 12 月にネット囲碁界で話題を席卷した出来事があった。謎の囲碁棋士 Master が突如現れてプロ棋士 60 人に短期間に連勝(60 勝 0 敗)したという[1]。筆者はリアルタイムでこの情報に接し、非常に大きな衝撃を受けた。同様のボードゲームであるチェスでは、1997 年に世界最強棋士カスパロフに IBM の開発したプログラムが勝利している。チェスと似ているが、取り上げた駒を再利用できる将棋は遥かに難しいといわれていたが、2006 年ごろから人間より強くなったといわれており、現在ではプロ棋士の研究に使われるほどだ(ちなみに AI で評価された先手角換わり戦法は小学生時代の得意技だった。えへん)。

将棋と比較すると、囲碁は手の候補が非常に多い。ルールをご存知でない読者のために簡単に説明すると、将棋の場合には今ある位置からの移動しか認められないが、囲碁の場合は盤面のどこに打つことも可能である。読む手の数が発散するので、コンピュータには 30 年くらいは負けない、といわれていた(年数は記憶があやしい)。また、将棋の場合には玉を詰めるという目的がわかりやすいが、囲碁の場合には短期的な目的が必ずしも勝利には結びつかないので、評価関数を作りづらいとされていた。

しかし、プロ棋士に 60 連勝である。囲碁界は当然騒然となった。正体は Google の作った囲碁 AI だということが明らかになった。その後、AlphaGo として世に現れ、2016 年 3 月に世界最強棋士李世ドルに勝利した。今やプロ棋士同士の対局でも参考にされている(NHK 番組など)。ちなみに、囲碁 AI を作ったハサビスはタンパク質の構造予測を正確にできる AlphaFold を世に生み出しおり、ノーベル賞候補との呼び声も高い[2]。

一連の出来事はガラス研究とは結びくと考えてはいなかった。一方、現職に異動した後、環境も変わったことだしこれまでと何か違った研究ができないだろうか、アイデアでリソース不足を補えるような研究はできないだろうかと模索していた。偶然、2016 年末ごろに京都大学物質-細胞統合システム拠点のバックウッド氏の研究[3](金属表面への有機分子の配向)についての講演に接した。囲碁 AI 誕生の記憶も新しかったことから、ガラスの物性予測に機械学習が使えるのではないかと考えたことが研究を始めたきっかけである。

2 機械学習とは？

機械学習とは、機械(コンピュータ)がデータから自動で学習し、その背景にあるパターンやルールを発見する方法というように説明される。コンピュータが自動で学習する、といっても、目の前の箱が何も指示を受けずに学習するわけではない。評価関数を定義し、評価値が最小(または最大)になるように関数を変えることによって学習する。

広義には、実験結果のフィッティングも機械学習に含まれる。例えば、組成 x に対して物性 y が一次関数 $y = ax + b$ に従うとしたときに、 a と b を決めるのも機械学習の一種である。このように考えると、読者の皆さんの研究でも機械学習を頻繁に使っているのではないだろうか。

3 データベースについて

機械学習においては大量のデータを必要とする。ガラスの分野においては 1930 年代よりデータ収集が行われ

ており、1980年代に O. V. Mazurin らによってハンドブックとして上市されている[4]。学生時代にこのハンドブックで何度も調べた記憶がある。INTERGLAD は、1991年に New Glass Forum のプロジェクトにより公開されたデータベースである[5]。このデータベースには、ガラスに関する論文や特許から組成や物性の情報が取り込まれているため、機械学習を行う環境として非常に良い。

4 データを用いた物性予測

ガラス分野ではデータベースを利用した物性予測が数多くなされている。例えば、2005年ごろに A. Flugel らがデータベースを利用した回帰による物性予測を行なっている[6]。我が国でも 2011年ごろに難波らのグループが INTERGLAD を用いた物性予測の先駆けとなる研究を行なっている[7]。この分野の進展は目覚ましく、ニューラルネットを複層重ねたディープラーニングを用いた物性予測についての報告も行われている[8]。

5 機械学習の基礎[9]

2でも述べたように、機械学習とは入力 x に対する y の予測である。 x と y の関係が図 1 のようになっている時、線形近似式 $y = w_0 + w_1x$ はエクセルを用いると容易に求めることができる ([10]よりダウンロード可能。例 1 参照)。 w_0 と w_1 が決まれば、 x に対する y の予測値を計算することができる。予測値と実データの誤差が最も小さくなるように w_0 と w_1 を決めることが線形近似で行っていることである。

ここで、 $\mathbf{X} = (\mathbf{1} \quad \mathbf{x})$ と定義する。ただし「 $\mathbf{1}$ 」は縦に数値 1 の並んだベクトルであり、 \mathbf{x} もデータを縦に並べたベクトルである。以降、太字でベクトルまたは行列を表す。次の式で w_0 と w_1 を求められることが知られている。

$$\mathbf{w} = \begin{pmatrix} w_0 \\ w_1 \end{pmatrix} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad \text{式(1)}$$

ここで T は転置、 $^{-1}$ は逆行列を示している ([9]例 2)。2 次の近似式 $y = w_0 + w_1x + w_2x^2$ の場合でも、 $\mathbf{X} = (\mathbf{1} \quad \mathbf{x} \quad \mathbf{x}^2)$ とすることで求めることができる ([9]例 3)。次数が上がると過学習が起きるので、それを防ぐための正則化が知られている。 λ を適切に設定することで、

$$\mathbf{w} = ((\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T) \mathbf{y} \quad \text{式(2)}$$

により \mathbf{w} が求まる。ただし \mathbf{I} は単位行列である ([10]例 4)。これをリッジ回帰という。

ここまで何を行ったかを復習するために、 \mathbf{x} をあるガラスの組成(百分率)、 \mathbf{y} をその物性としてみよう。ガラス組成と物性のデータの集まりがあったとすると、それらを \mathbf{x} 、 \mathbf{y} のように太字で表記する。2 次の近似の場合には、 $\mathbf{X} = (\mathbf{1} \quad \mathbf{x} \quad \mathbf{x}^2)$ である。 \mathbf{w} を係数として、組成 \mathbf{x} と物性 \mathbf{y} は

$$\mathbf{y} = \mathbf{X} \mathbf{w}^T \quad \text{式(3)}$$

の式で関係付けられている。 \mathbf{w} が式(1)によって決まると、物性 \mathbf{y} を予想することができる。過学習を防ぐには、 λ を適切に設定して式(2)により \mathbf{w} を求めれば良い。エクセルでも計算できるので、[9]例 2~4 を参照されたい。

6 カーネルリッジ回帰[11]

機械学習の中でも、統計的に解釈がしやすく、かつ非常に強力な手法の一つであるカーネルリッジ回帰につい

て平易に説明する。なお、カーネルリッジ回帰は、ガウス過程回帰と数学的に同等の手法である。

データ x_1 と x_2 の距離を

$$k_{1,2} = \exp(-\beta(x_1 - x_2)^2) \quad \text{式(5)}$$

と定義する。 x_1 と x_2 が等しければ 1, 離れるにしたがって急速に 0 となるので, 距離とみなせる。データ x_1, x_2, \dots, x_n に対してカーネル \mathbf{K} を次のように定義する。

$$\mathbf{K} = \begin{pmatrix} k_{1,1} & \cdots & k_{1,n} \\ \vdots & \ddots & \vdots \\ k_{n,1} & \cdots & k_{n,n} \end{pmatrix} \quad \text{式(6)}$$

カーネル回帰では, \mathbf{y} と \mathbf{x} の関係を (\mathbf{x} は \mathbf{K} を通じて) $\boldsymbol{\alpha}$ を回帰係数として

$$\mathbf{y} = \mathbf{K}\boldsymbol{\alpha}^T \quad \text{式(7)}$$

のように表すことができる。式(7)は, 線形回帰の式(3)と同じ形をしているので, 式(1)と同様に

$$\boldsymbol{\alpha} = (\mathbf{K}^T \mathbf{K})^{-1} \mathbf{K}^T \mathbf{y} \quad \text{式(8)}$$

として係数 $\boldsymbol{\alpha}$ を求めることができる。過学習を防ぎたい場合にも, やはり式(2)と同様に

$$\boldsymbol{\alpha} = (\mathbf{K}^T \mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{K}^T \mathbf{y} \quad \text{式(9)}$$

として係数 $\boldsymbol{\alpha}$ を求めることができる([9]例 5)。これがカーネルリッジ回帰である。

少し解説を加えておく。ガラスの組成を \mathbf{x}_1 から \mathbf{x}_n に対して物性が \mathbf{y}_1 から \mathbf{y}_n のように知られているとする。この時, 新たな組成 \mathbf{x}_m の物性は以下のように予測できる。

$$y_m = \alpha_{1,m} k_{1,m} + \alpha_{2,m} k_{2,m} + \cdots + \alpha_{n,m} k_{n,m} \quad \text{式(10)}$$

式(5)のように k は組成の距離を表すので, その距離に重み α を掛けて総和を取ることで予測が可能ということの意味している。先にも述べたように, $k_{n,m}$ は x_n と x_m が等しければ 1, 離れるにしたがって急速に 0 となるので, ガラス組成の似たものの影響が大きく, 似ていないものの影響が小さいことを反映させた予測である。

7 ガラス物性予測

カーネルリッジ回帰によって, 光学特性の一つであるアッベ数の予測を行った例について紹介する[12]。INTERGLAD で屈折率とアッベ数が報告されているガラスのうち, 屈折率が 1.9 から 2.0 かつ SiO_2 を組成に含むものとして 879 個のガラスデータを抽出した。ガラス組成の百分率を \mathbf{x} , アッベ数を \mathbf{y} としてカーネルリッジ回帰を行ったところ, 決定係数(1 に近いほど予測性能が高い)0.998 のモデルを作成することができた(図 2)。879 個のガラスに 56 組成の酸化物を 5%加えてデータベースには掲載されていないガラスを作製して実測し, 予測値との

比較を行うことにより、未知組成の物性予測が可能だということも明らかにした。

8 逆問題

ここまで組成から物性を予測するという話を解説してきた。これは比較的容易であり、順問題と呼ばれる。傾向を知ることにも一定の意味はあるが、データを用いた課題解決とはいえない。本来、我々が目指すべきは、目標とする物性があり、その物性を達成できる組成を明らかにすることである。これは逆問題と呼ばれている。

逆問題を解くアプローチの一つにベイズ最適化が知られている[3, 13]。例えば最大値を知りたい場合(アッベ数を大きくしたい場合など)、ガラスの組成空間は広いので、闇雲に探索することは好ましくない。予測値 y が大きくなる組成 x を探索するのも一つの方法であろう。もしくは、予測値 y は必ずしも大きくないが、探索が十分ではなく、突如として値が大きくなるかもしれない組成 x を探索するのも良い。これらを組み合わせて次の探索点を予測して実施することを繰り返して最大値に到達する。これがベイズ最適化である(図3)。

ここで、予測値としては上述したカーネルリッジ回帰の予測値を用いることができる。また、突如として大きくなるかもしれないという値は、分散(統計学でいう分散)を用いることで評価できる。分散は、ガウス過程回帰により得られる。これらをバランスよく組み合わせたものが期待改善度である。期待改善度の大きくなる x を順に試行することを繰り返すと、最大値に到達できる[14]。このアプローチにより、逆問題を解くことが可能である。

9 終わりに

データを用いたガラス材料探索の話をする、10人のうち2人くらいにはイヤな顔をされる。データで新しい材料は作れないよ、と。材料科学者としての経験を元にして新しいガラスを作るのが王道である。確かにそうであろう。しかし、場合によってはデータで傾向を知るだけで十分ということ“も”ある。両方のアプローチをうまく組み合わせることで、新しい材料が生まれることを期待している。

タイプライターからワードプロセッサ、そしてPCへと文章作成が変わっていったことを思い起こしてほしい。随分と便利になって、今やPCは必要不可欠な存在となった。また、電卓が現れた時にも、算盤の方が速いといわれたそうだ。機械学習も道具の一種である。活用できる道具はどんどん活用し、余剰のリソースを別の研究に注力する方が良い。囲碁AIや生成AIが広く使われているように、大量のデータを学習した予測モデルを誰もが容易に使える時代である。ましてや機械学習をノーコードで使えるのだから、これを利用しない手はないだろう。

参考文献

[1] <https://ja.wikipedia.org/wiki/AlphaGo>

[2] <https://alphafold.ebi.ac.uk>

[3] D.M. Packwood, *Bayesian Optimization for Materials Science*, (SpringerBriefs in the Mathematics of Materials) (Springer, 2017).

[4] O. V. Mazurin, M. V. Strelstina, and T. P. Shvaiko-Shvaikovskaya, *Handbook of Glass Data* (Elsevier, 1983)

[5] https://www.newglass.jp/interglad_n/gaiyo/info_j.html

[6] A. Flugel, *Glass Technol.: Eur. J. Glass Sci. Technol. A* **50** (2009)

[7] K. Ishii, T. Tsuneoka, S. Sakida, Y. Benino, and T. Nanba, *J. Ceram. Soc. Jpn.* **120**, 98 (2012)., 石井久美子, 恒岡徹, 崎田真一, 紅野安彦, 難波徳郎, *ニューガラス* **26** (2011)

[8] Y. Tokuda, M. Fujisawa, J. Ogawa, and Y. Ueda, *AIP Advances* **11**, 125127 (2021)

- [9]ビショップ C., パターン認識と機械学習 (丸善, 2012)., 持橋大地, 大羽成征, ガウス過程と機械学習 (講談社, 2014).
- [10] <https://glassl.net/sites/wp-content/uploads/2024/04/newglass2024.xlsx>
- [11]赤穂 昭太郎, カーネル多変量解析(岩波書店, 2008), 瀬戸 道生, 伊吹 竜也, 畑中 健志, 機械学習のための関数解析入門(内田老鶴圃, 2021)
- [12] Y. Tokuda, M. Fujisawa, D.M. Packwood, M. Kambayashi, and Y. Ueda, AIP Adv. **10**, 105110 (2020).
- [13]松井孝太, 金森研太, 豊浦和明, 竹内一郎, まてりあ **58**, 12 (2019)., マテリアルズ・インフォマティクス Q&A 集-解析実務と応用事例- (情報機構, 2020)
- [14]徳田陽明, 藤澤美沙, 武立奈々, 上田義勝, マテリアルステージ **23**(2023)

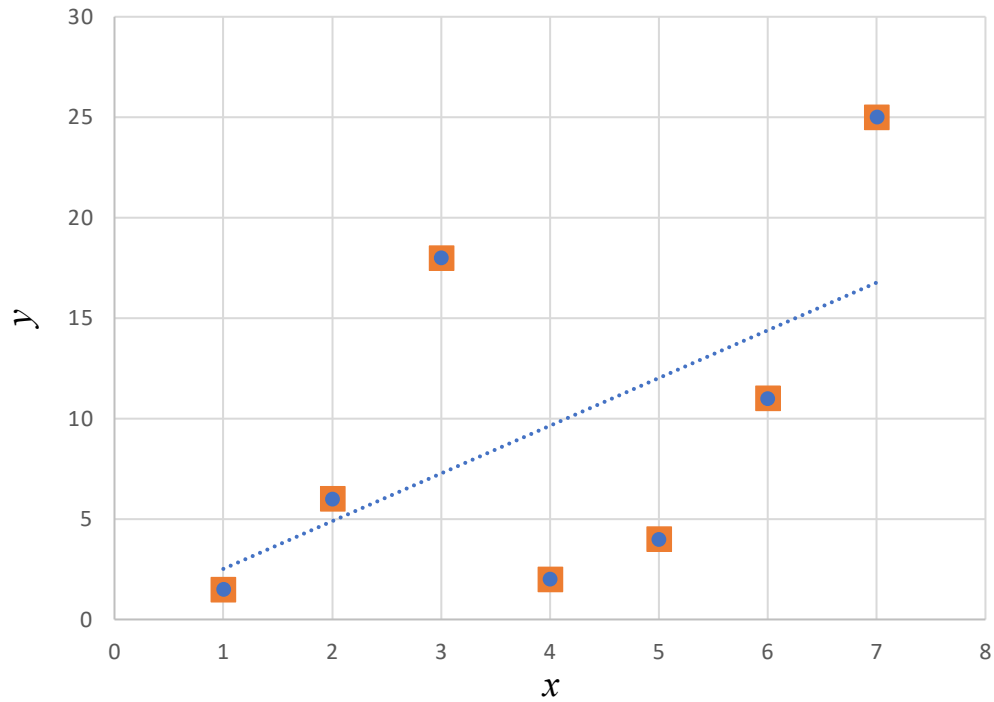


図1 x と y の散布図とそれらの線形近似。[8]からエクセルファイルをダウンロードできる。

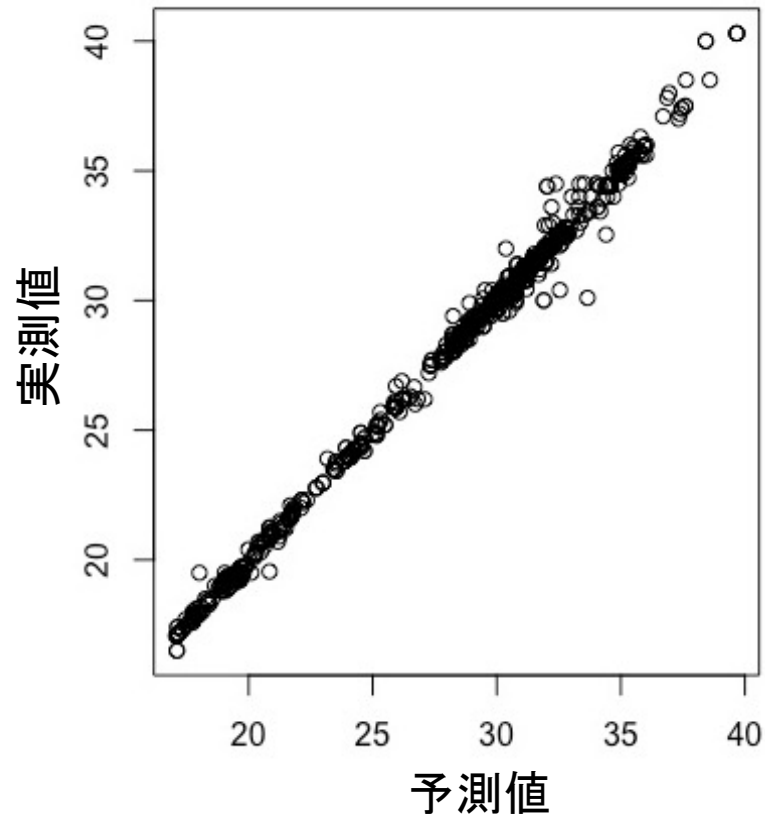


図2 カーネルリッジ回帰によって求めたアツベ数の予測値と実測値の比較。当てはまりの良さを示す R^2 値は0.998だった(1の時, 完全に一致)。

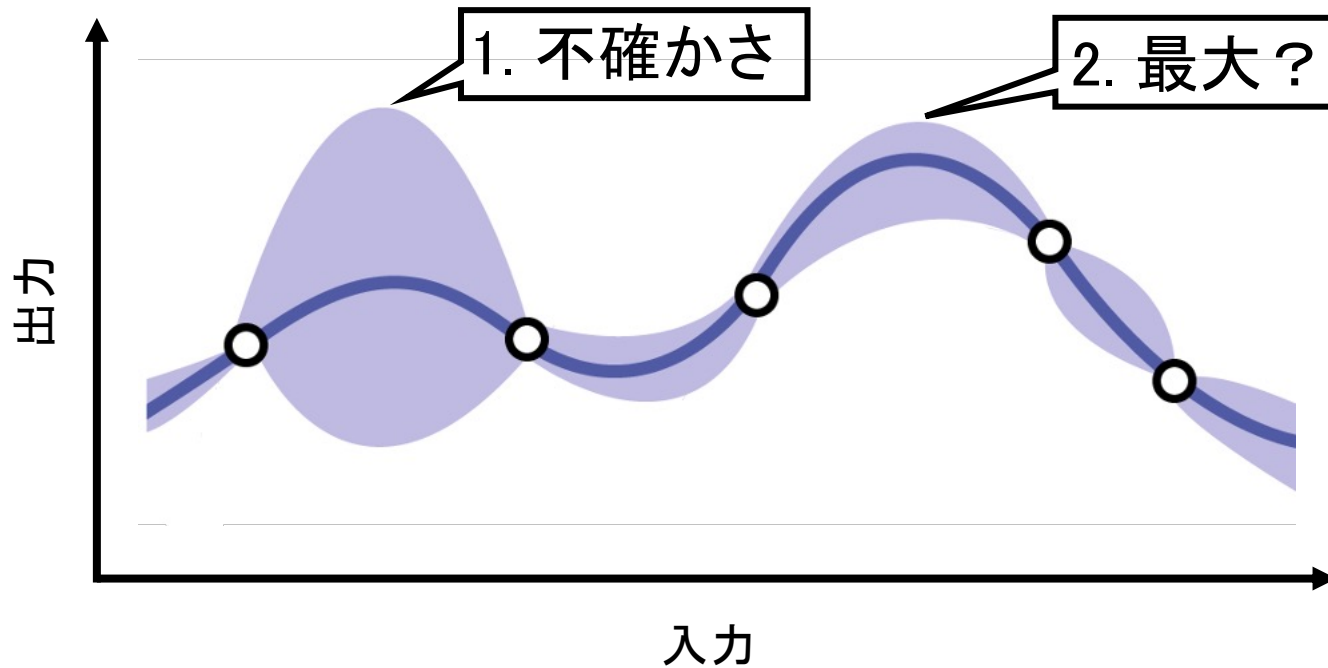


図3 説明変数が1つの場合にベイズ推定により最大値を得る例。太線が期待値であり、その上下に不確かさの度合い(分散)が示されている。○の点では組成と物性の関係が既に知られているので、不確かさが小さい。最大値となる物性を達成したい場合、1のような期待値が小さいが不確かさの大きな組成を探索するのか、2のような期待値が大きいが不確かさの小さな組成を探索するのかを期待改善度によって評価する。