

Three green apples are arranged in a cluster on a white background. One apple is in the foreground, slightly to the right, and two are behind it, one to the left and one to the right.

人工知能
IE229 - ARTIFICIAL INTELLIGENCE

第3回講義: 線形回帰
Lecture 3: Linear Regression

二宮 崇 (Takashi Ninomiya)
愛媛大学 (Ehime University)
ninomiya@cs.ehime-u.ac.jp

人工知能第3回の講義です。第3回は線形回帰について学びます。

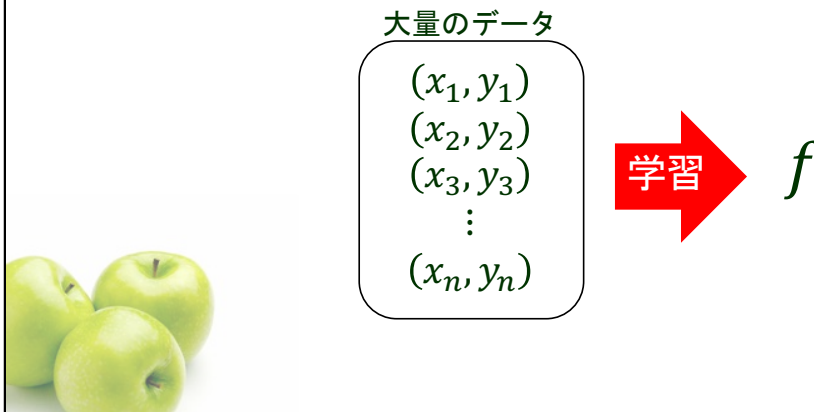
【復習】機械学習

- データから関数を学習

- 関数は、入力(x)と出力(y)の関係を表す

$$x \xrightarrow{f} y$$

- 大量の入出力ペア(x, y)の集まり(データ)から関数 f を予測



2

機械学習は、データから関数を学習することが本質です。関数は入力(x)と出力(y)の関係を表します。大量の入出力ペア(x, y)の集まり(データ)から入力と出力の関係をとらえた関数を自動的に獲得することが「学習」ということになります。

ただ単にデータがたくさんあれば良いわけではなく、この x のときはこんな y が正解、といったたくさんの入力と出力の例があることで、それらから、入力と出力の関係性を関数として学習する、というわけです。従って、学習には入力と出力のペアが揃った大量のデータが必要です。このように入力と出力のペアが揃ったデータから学習することを「教師つき学習」もしくは「教師あり学習」と言います。

【復習】機械学習

● データから学習

- 入力 $\mathbf{x} = (x_1, x_2, \dots, x_m)$ ← m 次元ベクトル
(画像の場合は、赤、青、緑の画像に対応する3つの行列(テンソル))

- データ D は入力 \mathbf{x} と出力 y のペアの集合

- 関数 f は **パラメータ(重み変数)の集合 $\mathbf{w} = (w_1, w_2, \dots, w_m)$** から成る

例: $f(\mathbf{x}) = w_1x_1 + w_2x_2 + \dots + w_mx_m$

- f は $f_{\mathbf{w}}(\mathbf{x})$ と書くとわかりやすい

- 学習 = データ D に対する誤差(損失)を最小にするパラメータ \mathbf{w} を求める

データ全体の誤差(損失)
$$L(\mathbf{w}) = \sum_{(\mathbf{x}, y) \in D} (y - f_{\mathbf{w}}(\mathbf{x}))^2$$



3

より一般的には、入力 \mathbf{x} はベクトルで与えられることが多いです。1次元、2次元ではなく、多次元のベクトルを与えます。一般に m 次元ベクトルとして与えます。画像の場合は、赤、緑、青の画像に対応する3つの行列(テンソル)として与えます。

データ D を入力 \mathbf{x} と出力 y のペアの集合とします。

関数 f はパラメータ(重み変数)の集合 ($\mathbf{w} = (w_1, w_2, \dots, w_m)$) からできていて、 \mathbf{w} と \mathbf{x} から y の値を予測することになります。例えば、 $f(\mathbf{x}) = w_1x_1 + w_2x_2 + \dots + w_mx_m$ といった形になります。 f は $f_{\mathbf{w}}(\mathbf{x})$ と書くとわかりやすいですね。

このように与えられたデータ D と関数 f に対し、データ D に対する誤差を最小にするパラメータ \mathbf{w} を求めることを、「学習」と言います。データ全体の誤差は、例えば、

$L(\mathbf{w}) = \sum_{(\mathbf{x}, y) \in D} (y - f_{\mathbf{w}}(\mathbf{x}))^2$ という式で計算できます。データの各ペア (\mathbf{x}, y) に対し、本当の正解の y から関数で予測される $f(\mathbf{x})$ の値を引くことで誤差が得られ、データ全体においてその二乗和を計算すると、それがデータ全体の誤差となります。この誤差のことを機械学習の分野では一般に「損失」と呼びます。データ全体に対する損失 $L(\mathbf{w})$ には \mathbf{x} と y と \mathbf{w} の3種類の変数が含まれますが、 \mathbf{x} と y にはデータから具体的な値が代入されますので、このデータ全体に対する損失は純粋にパラメータ \mathbf{w} だけの関数となります。これを最小化することでパラメータが求まります。

【復習】 回帰問題と分類問題と構造予測

関数: $y = f(x)$

- 回帰問題 (regression)

- $y \in R$ (実数)

例: 年齢予測、降水確率の予測、気温の予測

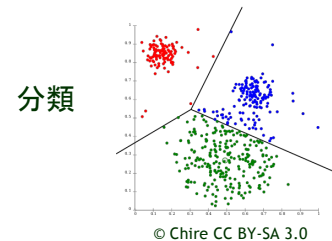
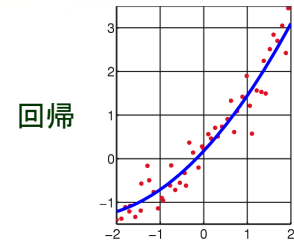
- 分類問題 (classification)

- $y \in \{C_1, C_2, \dots, C_K\}$ (ラベル集合)

例: 文書分類(政治、経済、スポーツ等)

- 構造予測 (structured prediction)

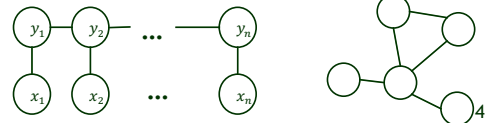
- $y \in G$ (グラフ集合)



© Chire CC BY-SA 3.0



構造予測



機械学習は関数 $y = f(x)$ を学習することでした。 x は様々な入力 that 想定されていますが、 y の値は大きく分けると3種類のタイプに分けられます。そのタイプによって、機械学習の問題は、回帰問題、分類問題、構造予測の3つに分けられます。

回帰問題は y が実数値になる場合です。例えば、年齢予測、降水確率の予測、気温の予測などがあげられます。右上の図のように x に対する y の値の曲線(関数)を学習することが目的となります。

分類問題は y がラベルになる場合です。「犬」「猫」ラベルのように、ラベル集合から正解を一つ選ぶ問題になります。ラベル集合は一般にこのように C_1, C_2, \dots, C_K として与えられて、この中から正解を一つ選ぶ、ということになります。例えば、文書分類問題があげられ、「政治」「経済」「スポーツ」などのラベルからその文書のジャンルを選択することになります。スライド右の真ん中にあるように、 x の空間をいくつかの領域に分けることが目的となります。

最後は構造予測と呼ばれる問題で、これは y の値がグラフになる場合です。グラフというのはこのスライドの図のように丸と直線で結ばれた化合物の模型のような形をした構造のことを指します。例えば、DNAの系列解析や機械翻訳などもこの構造予測の一種となります。図のように与えられた x に対応するグラフを生成することが目的となります。

損失関数

- 損失関数 $L(y, \hat{y})$

- 誤差を数値化する関数。 y は正解、 \hat{y} は予測値

- データ全体の損失 $L(\mathbf{w})$

$$L(\mathbf{w}) = \sum_{(x,y) \in D} L(y, f(x))$$

学習：データ全体の損失が最小になる各パラメータの値を探すこと

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} L(\mathbf{w})$$

- 回帰問題の損失関数(二乗損失)

$$L(y, \hat{y}) = (y - \hat{y})^2$$

- 分類問題の損失関数(負対数尤度、交差エントロピー損失)

$$L(y, \hat{y}) = -\log p_{\mathbf{w}}(y|\mathbf{x}) = \sum_{k=1}^K -y_k \log \hat{y}_k$$



5

続いて、学習の仕組みについて勉強していきます。学習はデータ全体の損失 $L(\mathbf{w})$ を最小化すること、と説明しました。データ全体の損失は、各データに対する損失を合計することで得られます。各データに対する損失を計算する関数 $L(y, \hat{y})$ のことを「損失関数」と呼びます(ただし、 y は正解、 \hat{y} は予測値)。問題の性質によっていろんな損失関数が用いられていて、例えば、回帰問題には先ほどまで説明した二乗損失が用いられ、分類問題では負対数尤度もしくは交差エントロピー損失と呼ばれる損失関数が用いられます。二乗損失は先ほどの説明と同じになりますが、正解の出力 y と関数が予測する $f(x)$ の値の差の二乗となります。データ D の入出力ペア (\mathbf{x}, y) に対して、これらの合計が最小になるように \mathbf{w} を調整すると考えれば良いです。

負対数尤度は分類問題に対して用いられる損失関数です。統計学における最尤推定を行うための損失関数となっています。データ全体の各データの尤度(確率)を負の対数にして最小化しています。対数をつけても求めるべきパラメータは同じになります。対数にするのは微分を簡単にするだけではなく、データ全体の尤度(確率)を計算機で表現するためでもあります。確率自体非常に小さな値であり、データ全体の確率は各データの確率の積となるため、計算機で表せられない非常に小さな値となってしまいます。対数をつけることで計算機で表現可能な値となります。例えば、 $\log_{10} 10^{-60} = -60$ となります。つまり、分類問題は対数尤度の最大化(最尤推定)になっていることに注意しましょう。

教師データに対する損失が最小になるように各パラメータの値を探すことが学習、ということになります。

今日の学習内容

- 回帰問題の解法を学びます
- 線形回帰モデル(入出力間の関係として線形関数を使った回帰モデル)を学びます
 - 線形回帰：最小二乗法、リッジ回帰、正規方程式



6

今回は回帰問題の解法について学びます。

具体的な例として線形回帰モデルを学びます。線形回帰モデルというのは、入出力間の関係として線形関数を用いた回帰モデルのことです。回帰モデルは回帰問題を解くためのモデルです。

線形回帰モデルを解くための最小二乗法、過学習を防ぐためのリッジ回帰、線形回帰モデルのための最小二乗法の解となる正規方程式について学びます。

回帰モデル: 学習

● 回帰モデルのための教師データの例 (学習時)

| 乗客数 | 離陸重量 | 燃料重量比 | 巡行燃費 | 航続距離 |
|-----|-------|-------|------|-------|
| 120 | 50.0 | 0.250 | 0.65 | 3850 |
| 108 | 60.6 | 0.309 | 0.6 | 5820 |
| 108 | 65.1 | 0.312 | 0.6 | 6610 |
| 496 | 375.0 | 0.414 | 0.55 | 11800 |

学習



回帰
モデル

入力
(特徴量ベクトル)

目標出力
(実数値)



7

まず、回帰問題のデータをみてみましょう。この表が回帰問題のデータとなっています。

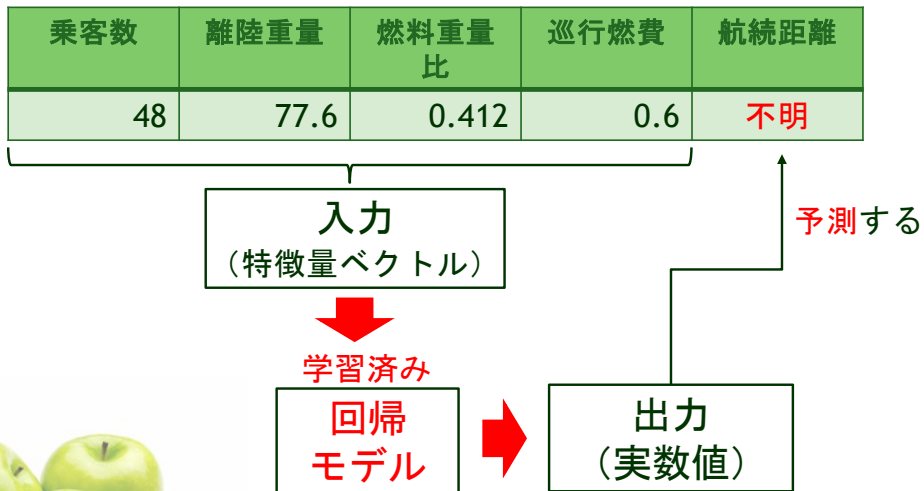
飛行機の航続距離に関するデータとなっていて、乗客数、離陸時の重量、燃料重量比、巡航燃費がわかっているときの実際の航続距離を表したデータとなっています。これらのデータのうち左4列の乗客数、離陸時の重量、燃料重量比、巡航燃費が回帰モデルへの入力データとなります。このとき、航続距離がいくらになるのか出力する回帰モデルを学習したいとします。

学習時は、これらの4つ特徴が入力として与えられ、航続距離が出力として与えられて学習を行います。この場合は4つのデータから成り立ちますが、これらのデータから回帰モデルを学習します。

回帰モデルの特徴として、出力が実数値となることに注意してください。(出力がラベルの場合は識別モデルと呼ばれます)

回帰モデル：推論

- 予測の例（推論時）



8

学習が終わった後、学習済みの回帰モデルを用いて予測を行います。未知のデータに対して、出力を予測することは「推論(inference)」と呼ばれます。

この例の場合、乗客数、離陸重量、燃料重量比、巡航燃費の4つの入力情報から成る未知のデータが与えられたとき、これらの入力を回帰モデルに与えて、出力値を得ます。この出力値が航続距離の予測値となります。

線形回帰モデル

- 入力ベクトル x から出力 y を得る関数が x の線形関数(w と x の内積)

$$y = w_0 + w_1x_1 + w_2x_2 + \cdots + w_Kx_K = \sum_{i=0}^K w_ix_i = \langle \mathbf{w}, \mathbf{x} \rangle = \mathbf{w}^T \mathbf{x}$$

$$\text{ただし、} \mathbf{x} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_K \end{pmatrix}, \mathbf{w} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_K \end{pmatrix}$$

T は転置

(w_0, w_1, \dots, w_K) のこと



※ x が1次元の時 : $y = w_1x_1 + w_0$

9

線形回帰モデルは具体的にはどんな形の関数でしょうか。

線形式が関数となっているモデルは線形モデルと呼ばれます。回帰問題の線形モデルは線形回帰モデルと呼ばれ、分類問題の線形モデルは線形識別モデルと呼ばれます。

線形式なので、入力を x_1, \dots, x_K としたとき、線形回帰モデルは $y = w_1x_1 + w_2x_2 + \cdots + w_Kx_K$ となります。 w_0, w_1, \dots, w_K は学習すべき重み変数(パラメータ)と呼ばれています。各 w_i は各入力 x_i に対する係数となっています。入力と重み変数をそれぞれ一つのベクトル

として表現すると、 $\mathbf{x} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_K \end{pmatrix}$ 、 $\mathbf{w} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_K \end{pmatrix}$ となります。機械学習の世界では一般的に

ベクトルは縦一列の行列として表現されることに注意してください。また、線形式の切片 w_0 を対応する入力として、 \mathbf{x} の最初の要素を1としている点にも注意してください。切片 w_0 のことは機械学習の世界ではバイアス項と呼ばれています。このような表現を与えたとき、線形式は \mathbf{w} と \mathbf{x} の内積 $\langle \mathbf{w}, \mathbf{x} \rangle$ になります。各ベクトルを縦一列の行列として表現すると、この内積の計算は $\mathbf{w}^T \mathbf{x}$ となります。線形式を簡潔に表現することができるため、 $\mathbf{w}^T \mathbf{x}$ と表現することが多いです。右肩の T は転置を表します。

x が1次元の時は $y = w_1x_1 + w_0$ となります。

線形回帰モデルの学習

- 学習: 損失の最小化

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} L(\mathbf{w})$$

$L(\mathbf{w})$: データ全体の損失

- 線形回帰モデルの損失関数

- 最小二乗法

$$L(\mathbf{w}) = \sum_{(\mathbf{x}, y) \in D} \underbrace{(y - \mathbf{w}^T \mathbf{x})^2}_{\substack{\text{正解} \\ \text{出力}} \quad \substack{\text{予測} \\ \text{出力}}}$$

- リッジ回帰

$$L(\mathbf{w}) = \sum_{(\mathbf{x}, y) \in D} (y - \mathbf{w}^T \mathbf{x})^2 + \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|^2}_{\text{正則化項}}$$

$$\|\mathbf{w}\|^2 = \left(\sqrt{w_0^2 + w_1^2 + \dots + w_K^2} \right)^2 = w_0^2 + w_1^2 + \dots + w_K^2$$



10

続いて、線形回帰モデルを教師データから学習することを考えます。

学習は、データ全体の誤差を最小化する重みベクトルを求めることで実現されます。正解出力と予測出力の差分のことを誤差または損失と呼びます。つまり、学習はデータ全体の損失の最小化として実現されます。

データ全体に対し、正解と予測の差の二乗和を最小化することで重み変数を求める手法は、最小二乗法と呼ばれます。線形回帰モデルを最小二乗法で求める際の損失は $L(\mathbf{w}) = \sum_{(\mathbf{x}, y) \in D} (y - \mathbf{w}^T \mathbf{x})^2$ として表現されます。つまり、データ D の各要素 (\mathbf{x}, y) に対して、正解出力である y から、線形回帰モデルでの予測値 $\mathbf{w}^T \mathbf{x}$ を引くと、これがデータ (\mathbf{x}, y) に対する損失となります。この損失をデータ全体で足し合わせると、 $\sum_{(\mathbf{x}, y) \in D} (y - \mathbf{w}^T \mathbf{x})^2$ となります。この式の中で \mathbf{x} と y は具体的な数値が入っているため、定数と考えてかまいません。従って、このデータ全体の損失を表す関数は残った変数である重み変数に対する関数となります。 $L(\mathbf{w})$ と表現されるように、重みベクトル \mathbf{w} に対する関数となっていることに注意してください。

線形回帰モデルの学習方法として、その他にリッジ回帰と呼ばれる学習方法があります。リッジ回帰の損失は、 $L(\mathbf{w}) = \sum_{(\mathbf{x}, y) \in D} (y - \mathbf{w}^T \mathbf{x})^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$ となります。最小二乗法

の損失に、 $\frac{\lambda}{2} \|\mathbf{w}\|^2$ という項がついています。 $\|\mathbf{w}\|^2 = \left(\sqrt{w_0^2 + w_1^2 + \dots + w_K^2} \right)^2 = w_0^2 +$

$w_1^2 + \dots + w_K^2$ となるため、各重み変数の二乗和となっています。これは過学習を抑えるための項で、正則化項と呼ばれたり、罰則項とも呼ばれています。過学習が起きるときは極端に重み変数の値が大きな値になることが知られているので、重み変数の値が大きくなりすぎたときには罰則として損失に加わるような仕組みになっています。 λ は、データに対するモデルの損失 $\sum_{(\mathbf{x}, y) \in D} (y - \mathbf{w}^T \mathbf{x})^2$ と正則化項 $\frac{\lambda}{2} \|\mathbf{w}\|^2$ のバランスを調整するパラメータで、人手で調整する必要があることから、ハイパーパラメータと呼ばれています。

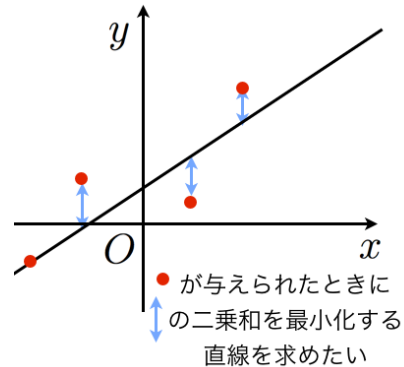
最小二乗法

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{(x,y) \in D} (y - \mathbf{w}^T \mathbf{x})^2$$

教師データDにおける入出力（ x と y ）の関係を
表すもっともらしい関数を求める

各教師データとの差の二乗和
を最小にする関数

例： x が1次元（モデルが直線）の時



11

最小二乗法は、従って、教師データDにおける入出力(x, y)の関係を表すもっともらしい関数を求めることになります。具体的には、モデルが表す関数の出力と各データの正解の差(図の水色矢印)が小さくなるように関数を調整することになります。

回帰モデルの学習例

- 教師データ $D = \{(1,2), (2,3), (3,5)\}$ から最小二乗法により学習される直線の式 $f(x) = ax + b$ を求めなさい。

$$L(a, b) = \sum_{(x,y) \in D} (y - f(x))^2$$



12

具体的に損失関数を最小化することで重み変数を求めてみましょう。損失関数に対する各重み変数の偏微分を全て0とすることで誤差関数を最小にする重み変数を求めます。

教師データ $D = \{(1,2), (2,3), (3,5)\}$ から最小二乗法により学習される直線の式 $f(x) = ax + b$ を求めましょう。(約5分)

回帰モデルの学習

- 教師データ $D = \{(1,2), (2,3), (3,5)\}$ から最小二乗法により学習される直線の式 $f(x) = ax + b$ を求めなさい。

$$L(a, b) = \sum_{(x,y) \in D} (y - f(x))^2$$

$$L(a, b) = (2 - a - b)^2 + (3 - 2a - b)^2 + (5 - 3a - b)^2$$

$$\begin{aligned} \frac{\partial L(a, b)}{\partial a} &= -2(2 - a - b) - 4(3 - 2a - b) - 6(5 - 3a - b) = \\ &= -46 + 28a + 12b = 0 \end{aligned}$$

$$\begin{aligned} \frac{\partial L(a, b)}{\partial b} &= -2(2 - a - b) - 2(3 - 2a - b) - 2(5 - 3a - b) = \\ &= -20 + 12a + 6b = 0 \end{aligned}$$

解くと、 $a = \frac{3}{2}, b = \frac{1}{3}$ 。従って、 $f(x) = \frac{3}{2}x + \frac{1}{3}$ 。

13

損失関数にデータの値を代入すると、 $L(a, b) = (2 - a - b)^2 + (3 - 2a - b)^2 + (5 - 3a - b)^2$ となります。

これを a と b でそれぞれ偏微分して 0 とする連立方程式を立てます。これを解くと、 $a = \frac{3}{2}, b = \frac{1}{3}$ となるため、求めるべき関数は $f(x) = \frac{3}{2}x + \frac{1}{3}$ となります。

ŵを求める...

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{(\mathbf{x}, \mathbf{y}) \in D} (\mathbf{y} - \mathbf{w}^T \mathbf{x})^2 = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{y} - X\mathbf{w}\|^2$$

$$\text{ただし、} \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}, \mathbf{w} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_K \end{pmatrix}, X = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1K} \\ \vdots & & \ddots & \vdots \\ 1 & x_{N1} & \cdots & x_{NK} \end{pmatrix}$$

$$\underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{y} - X\mathbf{w}\|^2 = \underset{\mathbf{w}}{\operatorname{argmin}} (\mathbf{y} - X\mathbf{w})^T (\mathbf{y} - X\mathbf{w})$$

$$\frac{\partial}{\partial \mathbf{w}} (\mathbf{y} - X\mathbf{w})^T (\mathbf{y} - X\mathbf{w}) = 0 \text{ を解けばよい}$$

$$\text{解くと... } (X^T X) \mathbf{w} = X^T \mathbf{y} \quad \text{正規方程式}$$

$$\hat{\mathbf{w}} = (X^T X)^{-1} X^T \mathbf{y}$$



14

線形回帰モデルには実は解析解が存在しています。誤差を各データごとに二乗して足し合わせていますが、これを行列形式で表現し、ベクトルで微分すると正規方程式と呼ばれる式が得られます。この正規方程式を解くことで解析解が得られます。

まず、重みベクトルは先程と同じ様に $\mathbf{w} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_K \end{pmatrix}$ とします。出力は全データの出力を

縦に並べたベクトル $\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}$ で表現します。入力各データを横ベクトル $\mathbf{x} =$

$(1 \ x_1 \ \cdots \ x_K)$ で表現し、全データの入力ベクトル \mathbf{x} を縦方向に並べることで行列にし

ます。つまり、 $X = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1K} \\ \vdots & & \ddots & \vdots \\ 1 & x_{N1} & \cdots & x_{NK} \end{pmatrix}$ となります。この行列はデータ行列と

呼ばれます。

そうすると、損失関数の最小化は、 $\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{(\mathbf{x}, \mathbf{y}) \in D} (\mathbf{y} - \mathbf{w}^T \mathbf{x})^2 = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{y} - X\mathbf{w}\|^2$ となります。 $\underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{y} - X\mathbf{w}\|^2 = \underset{\mathbf{w}}{\operatorname{argmin}} (\mathbf{y} - X\mathbf{w})^T (\mathbf{y} - X\mathbf{w})$ なので、

$\frac{\partial}{\partial \mathbf{w}} (\mathbf{y} - X\mathbf{w})^T (\mathbf{y} - X\mathbf{w}) = 0$ を解けば良いということになります。

この式を微分するのは難しいのですが、解くと、 $(X^T X) \mathbf{w} = X^T \mathbf{y}$ となります。この式は正規方程式と呼ばれています。この式に逆行列をかけることで、重みベクトルの解析解 $\hat{\mathbf{w}} = (X^T X)^{-1} X^T \mathbf{y}$ が得られます。

この正規方程式は難しい形をしていますが、元々、 $X\mathbf{w} = \mathbf{y}$ となる \mathbf{w} を見つけたかったわけなので、 $X\mathbf{w} = \mathbf{y}$ の左から X^T をかけて $(X^T X) \mathbf{w} = X^T \mathbf{y}$ としていると考えるとわかりやすいです。 $X\mathbf{w} = \mathbf{y}$ となる \mathbf{w} が直接見つければいいのですが、 X は正方行列ではないので、逆行列が存在しません。そのため、 X^T を左からかけることによって、無理矢理正方行列にしている、と考えるとわかりやすいです。つまり、 $(X^T X)^{-1} X^T$ は X^{-1} のような役割を果たしていると言えます。この擬似的な逆行列はムーア-ペンローズの擬似逆行列と呼ばれています。

【参考】正規方程式の導出

$$\frac{\partial}{\partial \mathbf{w}}(\mathbf{y} - X\mathbf{w})^T(\mathbf{y} - X\mathbf{w}) = 0$$

$$\frac{\partial}{\partial \mathbf{w}}(\mathbf{y}^T \mathbf{y} - \mathbf{w}^T X^T \mathbf{y} - \mathbf{y}^T X \mathbf{w} + \mathbf{w}^T X^T X \mathbf{w}) = 0$$

左辺の各項目の \mathbf{w} による偏微分

転置の公式: $(AB)^T = B^T A^T$

内積の微分の公式: $\frac{\partial}{\partial \mathbf{w}} \mathbf{w}^T \mathbf{a} = \frac{\partial}{\partial \mathbf{w}} \mathbf{a}^T \mathbf{w} = \mathbf{a}$ (\mathbf{a} は定数)

2次形式の微分の公式: $\frac{\partial}{\partial \mathbf{w}} \mathbf{w}^T A \mathbf{w} = (A + A^T) \mathbf{w}$

$$\frac{\partial}{\partial \mathbf{w}} \mathbf{y}^T \mathbf{y} = 0, \quad \frac{\partial}{\partial \mathbf{w}} \mathbf{w}^T X^T \mathbf{y} = X^T \mathbf{y}, \quad \frac{\partial}{\partial \mathbf{w}} \mathbf{y}^T X \mathbf{w} = [\mathbf{y}^T X]^T = X^T \mathbf{y}$$

$$\frac{\partial}{\partial \mathbf{w}} \mathbf{w}^T X^T X \mathbf{w} = (X^T X + (X^T X)^T) \mathbf{w} = 2X^T X \mathbf{w}$$

$$-2X^T \mathbf{y} + 2X^T X \mathbf{w} = 0 \text{ まとめると、} (X^T X) \mathbf{w} = X^T \mathbf{y}$$



15

参考までに正規方程式の導出を紹介しておきます。行列微分の公式を知らないといけないため、一般の大学生が習う知識で解くことは難しいです。

行列の公式として次のものを使います。

転置の公式: $(AB)^T = B^T A^T$

内積の微分の公式: $\frac{\partial}{\partial \mathbf{w}} \mathbf{w}^T \mathbf{a} = \frac{\partial}{\partial \mathbf{w}} \mathbf{a}^T \mathbf{w} = \mathbf{a}$ (\mathbf{a} は定数)

2次形式の微分の公式: $\frac{\partial}{\partial \mathbf{w}} \mathbf{w}^T A \mathbf{w} = (A + A^T) \mathbf{w}$

求めたい式は $\frac{\partial}{\partial \mathbf{w}}(\mathbf{y} - X\mathbf{w})^T(\mathbf{y} - X\mathbf{w}) = 0$ ですが、この式を展開して、 $\frac{\partial}{\partial \mathbf{w}}(\mathbf{y}^T \mathbf{y} - \mathbf{w}^T X^T \mathbf{y} - \mathbf{y}^T X \mathbf{w} + \mathbf{w}^T X^T X \mathbf{w}) = 0$ とします。

この式の各項に対して、ベクトル \mathbf{w} で微分をします。ちなみに、 $\frac{\partial f}{\partial \mathbf{w}} = \begin{pmatrix} \frac{\partial f}{\partial w_0} \\ \frac{\partial f}{\partial w_1} \\ \vdots \\ \frac{\partial f}{\partial w_K} \end{pmatrix}$ と定義され

ていますので、これも知らないといけないと解けません。

各項を公式を使って微分すると、 $\frac{\partial}{\partial \mathbf{w}} \mathbf{y}^T \mathbf{y} = 0$, $\frac{\partial}{\partial \mathbf{w}} \mathbf{w}^T X^T \mathbf{y} = X^T \mathbf{y}$, $\frac{\partial}{\partial \mathbf{w}} \mathbf{y}^T X \mathbf{w} = [\mathbf{y}^T X]^T = X^T \mathbf{y}$, $\frac{\partial}{\partial \mathbf{w}} \mathbf{w}^T X^T X \mathbf{w} = (X^T X + (X^T X)^T) \mathbf{w} = 2X^T X \mathbf{w}$ となり、まとめると、 $-2X^T \mathbf{y} + 2X^T X \mathbf{w} = 0$ となります。整理すると、 $(X^T X) \mathbf{w} = X^T \mathbf{y}$ になります。

x が1次元（直線）の場合

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}, \quad X = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix}$$

$$X^T X = \begin{bmatrix} N & \sum_{n=1}^N x_n \\ \sum_{n=1}^N x_n & \sum_{n=1}^N x_n^2 \end{bmatrix}, \quad (X^T X)^{-1} = \frac{1}{N \sum_{n=1}^N x_n^2 - (\sum_{n=1}^N x_n)^2} \begin{bmatrix} \sum_{n=1}^N x_n^2 & -\sum_{n=1}^N x_n \\ -\sum_{n=1}^N x_n & N \end{bmatrix},$$

$$X^T \mathbf{y} = \begin{bmatrix} \sum_{n=1}^N y_n \\ \sum_{n=1}^N x_n y_n \end{bmatrix}.$$

なので、 $\hat{w}_0 = \frac{1}{N} \left(\sum_{n=1}^N y_n - \hat{w}_1 \sum_{n=1}^N x_n \right),$

$$\hat{w}_1 = \frac{N \sum_{n=1}^N x_n y_n - \sum_{n=1}^N x_n \cdot \sum_{n=1}^N y_n}{N \sum_{n=1}^N x_n^2 - (\sum_{n=1}^N x_n)^2}.$$



16

x が1次元(直線)の場合は、 $X^T X$ が2x2の行列になるため、その逆行列を直接計算することができます。従って、解析解は

$$\hat{w}_0 = \frac{1}{N} \left(\sum_{n=1}^N y_n - \hat{w}_1 \sum_{n=1}^N x_n \right)$$

$$\hat{w}_1 = \frac{N \sum_{n=1}^N x_n y_n - \sum_{n=1}^N x_n \cdot \sum_{n=1}^N y_n}{N \sum_{n=1}^N x_n^2 - (\sum_{n=1}^N x_n)^2}$$

となります。

演習

- 教師データ= $\{(1,2), (2,3), (3,5)\}$ から最小二乗法により学習される直線の式を求めなさい。ただし、正規方程式の解を用いて答えなさい。



17

前に求めた最小二乗法の解を正規方程式の解を用いて解いてみましょう。
教師データ $D = \{(1,2), (2,3), (3,5)\}$ から最小二乗法により学習される直線の式 $f(x) = ax + b$ を求めましょう。(約5分)

演習（答え）

- 教師データ= {(1,2),(2,3),(3,5)} から最小二乗法により学習される直線の式を求めなさい。

$$N = 3, \quad \sum_{n=1}^N x_n = 6, \quad \sum_{n=1}^N y_n = 10, \quad \sum_{n=1}^N x_n y_n = 23, \quad \sum_{n=1}^N x_n^2 = 14$$

$$\hat{w}_1 = \frac{N \sum_{n=1}^N x_n y_n - \sum_{n=1}^N x_n \cdot \sum_{n=1}^N y_n}{N \sum_{n=1}^N x_n^2 - (\sum_{n=1}^N x_n)^2} = \frac{3 \times 23 - 6 \times 10}{3 \times 14 - 36} = \frac{9}{6} = \frac{3}{2}$$

$$\hat{w}_0 = \frac{1}{N} \left(\sum_{n=1}^N y_n - \hat{w}_1 \sum_{n=1}^N x_n \right) = \frac{1}{3} \left(10 - \frac{3}{2} \times 6 \right) = \frac{1}{3}$$



$$\text{答え} \quad \underline{\underline{\frac{3}{2}x + \frac{1}{3}}}$$

18

解は先程の解と同じになりますが、このようにして解析解を得ることができることがわかります。

リッジ回帰

- 最小二乗法に2次の正則化項を導入することで過学習を緩和

※過学習: 訓練データに過剰に適合してしまうこと

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \sum_{(x,y) \in D} (y - \mathbf{w}^T \mathbf{x})^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 \right\}$$

解くと



$$(X^T X + \lambda I) \mathbf{w} = X^T \mathbf{y}$$

$$\hat{\mathbf{w}} = (X^T X + \lambda I)^{-1} X^T \mathbf{y}$$

正則化項

$$\|\mathbf{w}\|^2 = \sqrt{w_0^2 + w_1^2 + \dots + w_K^2}$$

※ λ は正則化定数
(予め人手で定める)

※ I は単位行列



19

リッジ回帰に対する解析解も求めることができます。同様に解くと、正規方程式は

$$(X^T X + \lambda I) \mathbf{w} = X^T \mathbf{y}$$

となります。これに対し、逆行列を求めることで、

$$\hat{\mathbf{w}} = (X^T X + \lambda I)^{-1} X^T \mathbf{y}$$

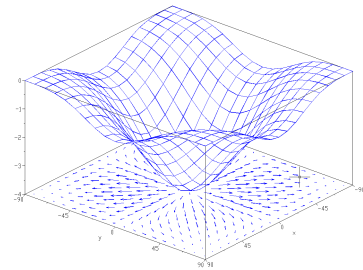
となります。ただし、 I は単位行列です。

最適化に関して

- 逆行列が簡単に求まる場合: 正規方程式を利用
- 逆行列演算が難しい場合: 数値計算法を利用
 - ニュートン法、準ニュートン法
 - 最急降下法
 - オンライン学習 (確率的勾配法(SGD), AdaGrad, Adamなど)

- 勾配を用いて計算する手法が多い

$$\text{勾配} \quad \nabla L = \left(\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \dots, \frac{\partial L}{\partial w_m} \right)$$



© Simiprof CC 表示 3.0

20

逆行列が簡単に求まる場合は、正規方程式を用いて解くときれいな解を高速に求めることができます。しかし、一般には逆行列を求めることが難しいため、数値計算法を利用して解を求めることになります。

数値計算法としては、ニュートン法、準ニュートン法、最急降下法が有名です。また、これらの学習方法はバッチ学習と呼ばれるデータ全体に対して最適化する手法が用いられていますが、深層学習では、一部のデータに対して最適化を繰り返すオンライン学習が一般的によく用いられています。オンライン学習では、確率的勾配降下法(SGD)、AdaGrad、Adamなどがあります。これらの手法はいずれも勾配に基づいて損失関数の最小化を行います。

まとめ

- **線形回帰モデル**

- 最小二乗法
- リッジ回帰
- 正規方程式

- **最適化**

- 正規方程式による最適解の導出
- 逆行列演算が難しい場合は勾配に基づく数値計算法



21

今回のまとめです。

今回は線形回帰モデルについて勉強しました。最小二乗法とリッジ回帰について学んで、その解析解として正規方程式を学びました。最適化には正規方程式を解くことで最適解を導出する方法がありますが、逆行列の計算が必要です。逆行列演算が難しい場合は勾配に基づく数値計算法でその解を求めます。