

機械学習・要点まとめ

第1章 線形回帰モデル

□ 線形回帰モデル

以下の線形結合において、最小二乗法によりパラメータ w_i を推定

$$\hat{y} = w^T x + w_0 = \sum_{j=1}^m w_j x_j + w_0 \quad (m: \text{パラメータ数、説明変数の数})$$

$$\text{回帰係数 } \hat{w} = (X^{(train)T} X^{(train)})^{-1} X^{(train)T} y^{(train)}$$

$$\text{予測値 } \hat{y} = X(X^{(train)T} X^{(train)})^{-1} X^{(train)T} y^{(train)}$$

$n \times (m+1) \quad (n: \text{データ数})$

□ 線形回帰モデルのハンズオン

- sklearn の場合: linear_model の LinearRegression を用いる。fit 関数でパラメータを推定、predict 関数で予測。 サマリー・考察は、skl_regression.ipynb に記載
- numpy の場合: $\hat{w}_1 = \text{Cov}[x, y] / \text{Var}[x]$, $\hat{w}_0 = \mu_y - \hat{w}_1 \mu_x$ (Cov: 共分散、Var: 分散、 μ : 平均) により、パラメータを推定する実装を行う。

第2章 非線形回帰モデル

□ 非線形回帰モデル

- 基底展開法: 回帰関数として、基底関数と呼ばれる既知の非線形関数とパラメータベクトルの線形結合で表現。未知パラメータは、線形回帰モデルと同様に最小二乗法や最尤法により推定。

$$\hat{y}_i = w^T (x_i) + w_0 = \sum_{j=1}^m w_j \varphi_j(x_i) + w_0$$

$$\text{説明変数} \quad x_i = (x_{i1}, x_{i2}, \dots, x_{im}) \in \mathbb{R}^m$$

$$\text{非線形関数ベクトル} \quad (x_i) = (\varphi_1(x_i), \varphi_2(x_i), \dots, \varphi_k(x_i)) \in \mathbb{R}^k$$

$$\text{非線形関数の計画行列} \quad \Phi^{(train)} = ((x_1), (x_2), \dots, (x_n)) \in \mathbb{R}^{n \times k}$$

$$\text{最尤法による予測} \quad \hat{y} = \Phi (\Phi^{(train)T} \Phi^{(train)})^{-1} \Phi^{(train)T} y^{(train)}$$

- よく使われる基底関数:

$$\text{多項式関数 } \varphi_j = x^j \quad \text{ガウス型基底関数 } \varphi_j(x) = \exp\left\{-(x - \mu_j)^2 / 2h_j\right\}$$

- 不要な基底関数を削除: 基底関数の数、位置やバンド幅によりモデルの複雑さが変化。多くの基底関数を用意してしまうと過学習が起こるため、適切な基底関数を用意(交差検証 CV 等で選択)。
- 正則化法(罰則化法): 「モデルの複雑さに伴って、その値が大きくなる正則化項(罰則項)を課した関数」を最小化。

$$S_y = (y - \Phi w)^T (y - \Phi w) + \gamma R(w) \quad (\gamma > 0)$$

基底関数の数(k)が増加するとパラメータが増加し、残差は減少(モデルが複雑化) モデルの複雑さに伴う罰則

□ 非線形回帰モデルのハンズオン

- サマリー・考察は、skl_nonlinear_regression.ipynb に記載

第3章 ロジスティック回帰モデル

ロジスティック回帰モデル

- 分類問題(クラス分類): **対数オッズ** $\log(p/1-p)$ を線形回帰で予測し、それを**正規化**して確率として出力を得ることで、結果としてクラス分類を実現。対数オッズの逆関数が**シグモイド関数**(活性化関数) 活性化関数(= **ロジット変換**)により、クラス i に属する確率 p が求まる。
- 目的関数は、**対数尤度**(回帰モデル構築時に最大化したい**尤度関数の対数**をとったもの、その負値が**交差エントロピー関数**)。交差エントロピー関数の最小化は、**対数尤度の最大化**と同値。

説明変数 $\mathbf{x} = (x_1, x_2, \dots, x_m)^T \in \mathbb{R}^m$

目的変数 $y \in \{0, 1\}$

パラメータ $\mathbf{w} = (w_1, w_2, \dots, w_m)^T \in \mathbb{R}^m$

線形結合 $z = \mathbf{w}^T \mathbf{x} + w_0 = \sum_{j=1}^m w_j x_j + w_0$

出力 $\hat{y} = P(Y=1 | \mathbf{x}) = \sigma(z = \sum_{j=1}^m w_j x_j + w_0) = \frac{1}{1 + \exp(-az)}$ σ : シグモイド関数

シグモイド関数の出力を $Y=1$ になる確率に対応させる

(Y は、 $\hat{y} \geq 0.5$ なら 1、 $\hat{y} < 0.5$ なら 0 と予測)

- ロジスティック回帰モデルの最尤推定: 対数尤度関数 E を最大とするパラメータを探索

$y_1 \sim y_n$ のデータが得られた際の尤度関数 $L(\mathbf{w})$ (確率 P はベルヌーイ分布に従うものとする):

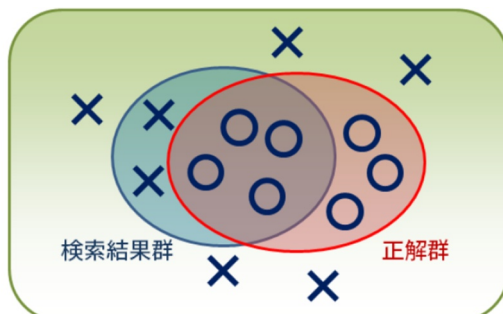
$$L(\mathbf{w}) = P(y_1, y_2, \dots, y_n | w_0, w_1, \dots, w_m) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i} = \prod_{i=1}^n \sigma(w^T x_i)^{y_i} (1 - \sigma(w^T x_i))^{1-y_i}$$

$$\text{交差エントロピー } E(\mathbf{w}) = -\log L(\mathbf{w}) = -\sum_{i=1}^n \{y_i \log p_i + (1 - y_i) \log(1 - p_i)\}$$

$$\frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} = -\sum_{i=1}^n (y_i - p_i) x_i \text{ より、}$$

$$\text{勾配降下法 } \mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \eta \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} = \mathbf{w}^{(k)} + \eta \sum_{i=1}^n (y_i - p_i) x_i$$

適合率と再現率



混同行列: 真陽性(TP)、偽陽性(FP)、偽陰性(FN)、真陰性(TN)

正解率 accuracy = (TP + TN) / 全体 = 9 / 14 = 0.64

適合率 precision = TP / (TP + FP) 偽陽性(FP)

再現率 recall = TP / (TP + FN) 偽陰性(FN)

$$\text{適合率 Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{4}{6} = 0.67 \quad \text{再現率 Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{4}{7} = 0.57$$

ロジスティック回帰モデルのハンズオン

- サマリー・考察は、skl_logistic_regression.ipynb , np_logistic_regression.ipynb に記載

機械学習・要点まとめ

第4章 主成分分析

□ 主成分分析

- 主成分分析 (PCA; Principal Component Analysis) : 教師なし学習で、多変量データの持つ構造をより少数個の指標に圧縮しつつそれに伴う情報の損失はなるべく小さくするため、変換後の分散を最大にするように探索する。

学習データ $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{im}) \in \mathbb{R}^m$

平均(ベクトル) $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$

データ行列 $\bar{\mathbf{X}} = (\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}})^T$

分散共分散行列 $= \text{Var}(\bar{\mathbf{X}}) = \frac{1}{n} \bar{\mathbf{X}}^T \bar{\mathbf{X}}$

線形変換後のベクトル $\mathbf{s}_j = (s_{1j}, s_{2j}, \dots, s_{nj})^T = \bar{\mathbf{X}} \mathbf{a}_j$ ($\bar{\mathbf{X}}: n \times m, \mathbf{a}_j: m \times 1$) j は射影後のインデックス

$\text{Var}(\mathbf{s}_j) = \frac{1}{n} \mathbf{s}_j^T \mathbf{s}_j = \frac{1}{n} (\bar{\mathbf{X}} \mathbf{a}_j)^T (\bar{\mathbf{X}} \mathbf{a}_j) = \frac{1}{n} \mathbf{a}_j^T \bar{\mathbf{X}}^T \bar{\mathbf{X}} \mathbf{a}_j = \mathbf{a}_j^T \text{Var}(\bar{\mathbf{X}}) \mathbf{a}_j$

線形変換後の変数の分散が最大となる射影軸は、以下のラグランジュ関数を最大にする:

$E(\mathbf{a}_j) = \mathbf{a}_j^T \text{Var}(\bar{\mathbf{X}}) \mathbf{a}_j - \lambda(\mathbf{a}_j^T \mathbf{a}_j - 1)$

$\frac{\partial E(\mathbf{a}_j)}{\partial \mathbf{a}_j} = 2 \text{Var}(\bar{\mathbf{X}}) \mathbf{a}_j - 2\lambda \mathbf{a}_j = 0 \quad \text{Var}(\bar{\mathbf{X}}) \mathbf{a}_j = \lambda \mathbf{a}_j$

$\text{Var}(\mathbf{s}_j) = \mathbf{a}_j^T \text{Var}(\bar{\mathbf{X}}) \mathbf{a}_j = \lambda_j \mathbf{a}_j^T \mathbf{a}_j = \lambda_j$ 射影先の分散は固有値と一致

- 寄与率: 各成分の重要度が測れる

第1～元次元分の主成分の分散は固有値と等しく、元のデータの分散と一致

元データの総分散 $V_{\text{total}} = \sum_{i=1}^m \lambda_i$

寄与率 c_k : 第 k 主成分の分散の全分散に対する割合 (第 k 主成分が持つ情報量の割合)

累積寄与率 r_k : 第1～ k 主成分まで圧縮した際の情報損失量の割合

$$c_k = \frac{\lambda_k}{\sum_{i=1}^m \lambda_i} \quad r_k = \frac{\sum_{j=1}^k \lambda_j}{\sum_{i=1}^m \lambda_i}$$

- オートエンコーダと主成分分析の関係・・・オートエンコーダの中間層の活性化関数を恒等関数にすると、オートエンコーダと主成分分析は等価な数式で表すことができ同じ結果を返す。

□ 主成分分析のハンズオン

- サマリー・考察は、`skl_pca.ipynb`, `np_pca.ipynb` に記載

第5章 k 近傍法

□ k 近傍法

- 分類問題で適用: 最近傍のデータを k 個取ってきて、それらが最も多く所属するクラスに識別
 k (ハイパーパラメータ) により結果が変わる、 k を大きくすると決定境界は滑らかになる

□ k 近傍法のハンズオン

機械学習・要点まとめ

- ・ サマリー・考察は、np_knn.ipynb に記載

第6章 k-means (k 平均クラスタリング)

▣ k-means

- ・ 教師なし学習、分類問題で適用： 最近傍のデータを k 個取ってきて、それらが最も多く所属するクラスに識別 k(ハイパーパラメータ)により結果が変わる、k を大きくすると決定境界は滑らかになる

【学習】

各クラスタ中心の初期値を設定する。

各データ点に対して、各クラスタ中心との距離を計算し、最も距離が近いクラスタを割り当てる

各クラスタの平均ベクトル(中心)を計算する

収束するまで ・ の処理を繰り返す

▣ k-means のハンズオン

- ・ サマリー・考察は、skl_kmeans.ipynb , np_kmeans.ipynb に記載

第7章 サポートベクターマシン

▣ サポートベクターマシン(線形・非線形共通)

【学習】

特徴空間上で線形なモデル $y(x) = \mathbf{w}^T \phi(x) + b$ を用い、その正負によって 2 値分類を行うことを考える。
サポートベクターマシンではマージンの最大化を行うが、それは以下の最適化問題を解くことと同じ。
ただし、訓練データを $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^T$, $\mathbf{t} = [t_1, t_2, \dots, t_n]^T$ ($t_i = \{-1, +1\}$) とする。

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2, \quad \text{subject to} \quad t_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1 \quad (i=1, 2, \dots, n)$$

< 線形の場合は、 $\phi(\mathbf{x}_i) = \mathbf{x}_i$ で、以降、同様に読み替えること >

ラグランジュ乗数法を使うと、上の最適化問題はラグランジュ乗数 \mathbf{a} (≥ 0) を用いて、以下の目的関数を最小化する問題となる。

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n a_i t_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b - 1) \quad \dots (1)$$

目的関数が最小となるのは、 \mathbf{w}, b に関して偏微分した値が 0 となるときなので、

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^n a_i t_i \phi(\mathbf{x}_i) = 0, \quad \frac{\partial L}{\partial b} = \sum_{i=1}^n a_i t_i = 0$$

これを式(1) に代入することで、最適化問題は結局以下の目的関数の最大化となる。

$$\tilde{L}(\mathbf{a}) = \sum_{i=1}^n a_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n a_i a_j t_i t_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = \mathbf{a}^T \mathbf{1} - \frac{1}{2} \mathbf{a}^T H \mathbf{a} \quad (1: \text{全成分が 1 のベクトル})$$

ただし、行列 H の i 行 j 列成分は $H_{ij} = t_i t_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = t_i t_j k(\mathbf{x}_i, \mathbf{x}_j)$ である。

また制約条件は、 $\mathbf{a}^T \mathbf{t} = 0$ ($\frac{1}{2} \|\mathbf{a}^T \mathbf{t}\|^2 = 0$) である。

機械学習・要点まとめ

この最適化問題を最急降下法で解く。目的関数と制約条件を a で微分すると、

$$\frac{d\tilde{L}}{da} = 1 - Ha, \quad \frac{d}{da} \left(\frac{1}{2} \|a^T t\|^2 \right) = (a^T t) t$$

従って、 a を以下の二式で更新する。

$$a \leftarrow a + \eta_1 (1 - Ha), \quad a \leftarrow a - \eta_2 (a^T t) t$$

【予測】

新しいデータ点 x に対しては、 $y(x) = w^T \phi(x) + b = \sum_{i=1}^n a_i t_i k(x, x_i) + b$ の正負によって分類する。

ここで、最適化の結果得られた $a_i (i=1, 2, \dots, n)$ の中で $a_i=0$ に対応するデータ点は予測に影響を与えないので、 $a_i>0$ に対応するデータ点 (**サポートベクトル**) のみ保持しておく。 b はサポートベクトルのインデックスの集合を S とすると、 $b = \frac{1}{S} \sum_{s \in S} (t_s - \sum_{i=1}^n a_i t_i k(x, x_i))$ によって求める。

□ ソフトマージン SVM

【学習】

分離不可能な場合は学習できないが、データ点がマージン内部に入ることや誤分類を許容することでその問題を回避する。

スラック変数 $\xi_i \geq 0$ を導入し、マージン内部に入った点や誤分類された点に対しては、 $\xi_i = |1 - t_i y(x_i)|$ とし、これらを許容する代わりに、ペナルティを与えるように、最適化問題を以下のように修正する。

$$\min_{w, b} \quad \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

$$\text{subject to} \quad t_i (w^T \phi(x_i) + b) \geq 1 - \xi_i \quad (i=1, 2, \dots, n)$$

ただし、パラメータ C はマージンの大きさと誤差の許容度のトレードオフを決めるパラメータである。

この最適化問題にラグランジュ乗数法などを用いると、結局最大化する目的関数はハードマージン SVM と同じになる。

$$\tilde{L}(a) = \sum_{i=1}^n a_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n a_i a_j t_i t_j \phi(x_i)^T \phi(x_j)$$

ただし、制約条件が $a_i \geq 0$ の代わりに $0 \leq a_i \leq C$ ($i=1, 2, \dots, n$) となる。(ハードマージン SVM と同じ $\sum_{i=1}^n a_i t_i = 0$ も制約条件)

- C について: C が大きいと完全に分離することに過度になり、汎化性能が小さくなるので、一般に過学習を防ぐためには C を小さくする方がよい。

□ サポートベクターマシンのハンズオン

- サマリー・考察は、np_svm.ipynb に記載

機械学習・要点まとめ

<http://study-ai.com/jdla/>



3ヶ月で現場で潰しが効く
ディープラーニング講座

 Study-AI

詳しくはこちら ▶