

クレジット:

UTokyo Online Education データマイニング入門 2018 森 純一郎

ライセンス:

利用者は、本講義資料を、教育的な目的に限ってページ単位で利用することができます。特に記載のない限り、本講義資料はページ単位でクリエイティブ・コモンズ 表示-非営利-改変禁止 ライセンスの下に提供されています。

<http://creativecommons.org/licenses/by-nc-nd/4.0/>

本講義資料内には、東京大学が第三者より許諾を得て利用している画像等や、各種ライセンスによって提供されている画像等が含まれています。個々の画像等を本講義資料から切り離して利用することはできません。個々の画像等の利用については、それぞれの権利者の定めるところに従ってください。



データマイニング入門 第9回

2018年度

学習目標

機械学習の概念について理解する

教師あり学習・回帰の概念について理解する

線形回帰

- コスト関数の最小化について理解する
- パラメータの推定
 - 最急降下法について理解する
- パラメータの解析解
 - 正規方程式について理解する

Pythonで線形回帰の実装を理解する

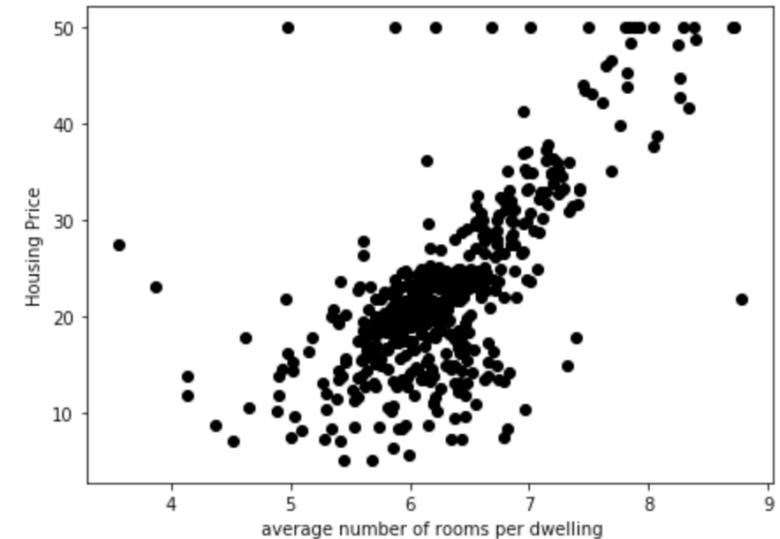
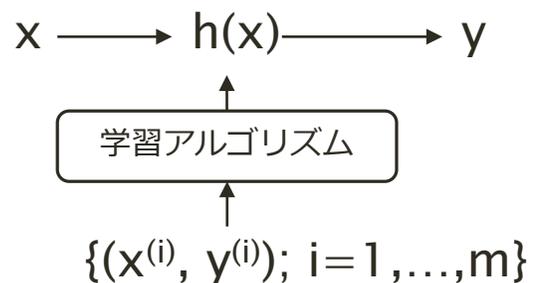
教師あり学習

入力から出力を予測したい

- 特徴量ベクトル（入力）： $x^{(i)}$
- 正解（出力）： $y^{(i)}$
- 訓練データ： $(x^{(i)}, y^{(i)})$
- （訓練）データセット： $\{(x^{(i)}, y^{(i)}); i=1, \dots, m\}$
 - i は学習データセットのインデックス

教師あり学習（supervised learning）

- データセットを元に入力から出力を予測する関数 $h(x)$ を学習する
 - h は仮説関数と呼ばれる



教師あり学習

特徴量ベクトル（入力）： $x^{(i)}$

正解（出力）： $y^{(i)}$

- ・ 連続値：回帰
- ・ 離散値：分類

データセット： $\{(x^{(i)}, y^{(i)}); i=1, \dots, m\}$

仮説関数 $h(x)$ のパラメータ： θ

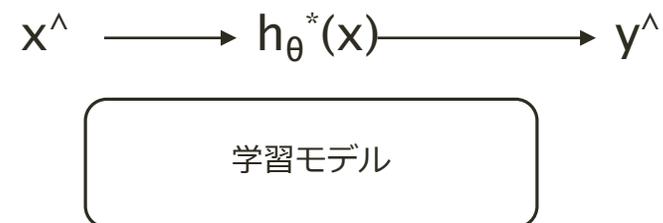
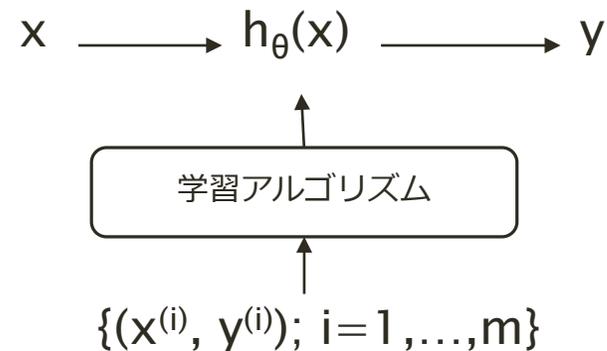
学習アルゴリズム

- ・ データセットを元に仮説関数 $h(x)$ のパラメータを学習
 - ・ コスト関数に基づくパラメータの最適化

未知の入力： x^\wedge 予測： y^\wedge

学習モデル

- ・ 未知の入力に対して尤もらしい出力を予測
- ・ $h(x)$ （最適パラメータ θ^* ）を元に入力 x^\wedge から y^\wedge を予測

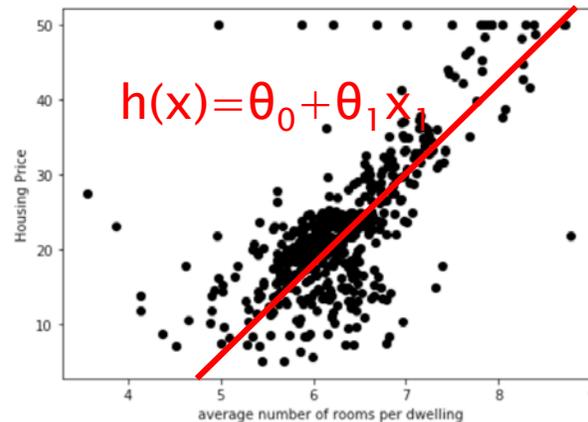


線形回帰

1変数の入力に対して出力を予測する線形回帰を考える

仮説関数 $h(x)$ の入力に対するパラメータ（重み）を θ_i として $h(x)$ は

- $h(x) = \theta_0 + \theta_1 x_1 = \sum \theta_k x_k = \theta^T x$
 - $\theta = (\theta_0, \theta_1)^T$
 - $x = (x_0, x_1)^T$ (バイアス項 $x_0 = 1$)
- $y = h(x) = b + ax$ と同じ
 - θ_0 は切片 b , θ_1 は傾き a



データセットを元に入力 x から出力 y を予測する関数 $h(x)$ のパラメータ θ を学習
学習した $h_{\theta}^*(x)$ を元に新たな入力 x^{\wedge} に対する出力 y^{\wedge} を予測したい

線形回帰

コスト関数

どのように仮説関数のパラメータを学習するか

- データセット $\{(x^{(i)}, y^{(i)}); i=1, \dots, m\}$ について
- 仮説関数の出力 $h(x^{(i)})$ と実際の出力 $y^{(i)}$ になるべく近くなるようなパラメータ θ を見つけたい

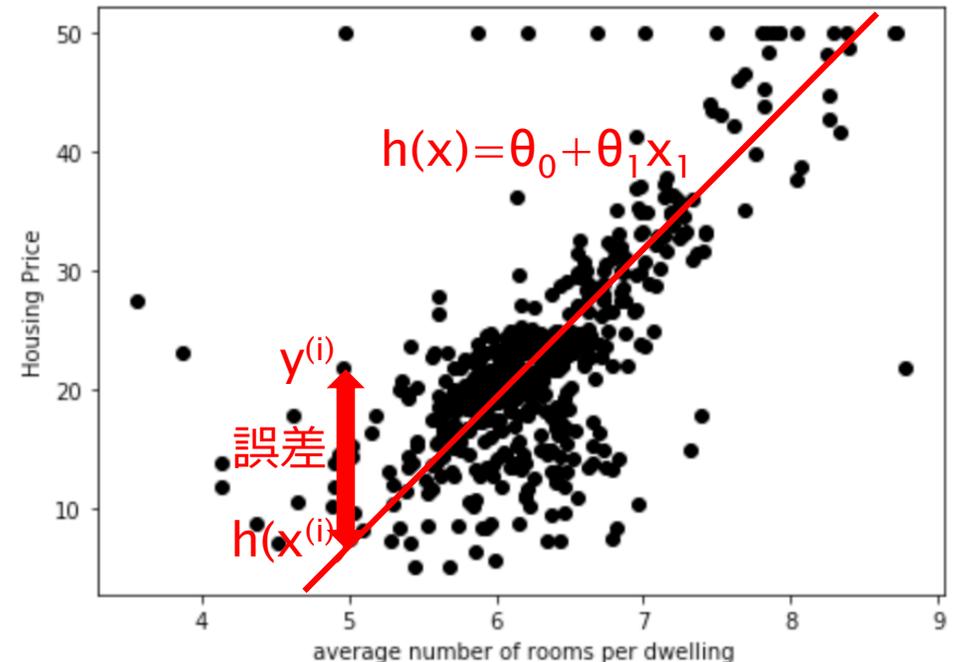
最小二乗法

- コスト関数（残差平方和）の最小化

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$$

データセット全体について $h(x^{(i)})$ と $y^{(i)}$ の誤差を少なくするようなパラメータ θ を見つけたい

入力が1変数の場合は誤差を少なくするような直線のパラメータ、傾き θ_1 と切片 θ_0 、を見つかることになる



線形回帰

コスト関数

入力の変数が多変数（ n 次元の特徴量ベクトル）の時、仮説関数 $h(x)$ は

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

この時、コスト関数は入力が1変数の時と同様に

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$$

m 個の $n+1$ 次元入力データを以下の様に表すと（ x_0 はバイアス項で $x_0=1$ とする）

$$X = \begin{pmatrix} x_0^{(1)} & x_1^{(1)} & \dots & x_n^{(1)} \\ \dots & \dots & \dots & \dots \\ x_0^{(m)} & x_1^{(m)} & \dots & x_n^{(m)} \end{pmatrix}$$

コスト関数は

$$J(\theta) = \frac{1}{2m} (X\theta - y)^T (X\theta - y)$$

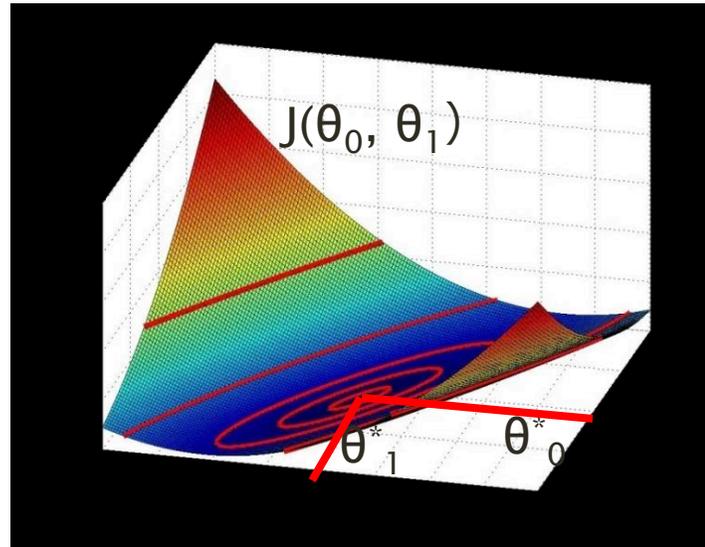
$$\theta = (\theta_0, \theta_1, \dots, \theta_n)^T$$

$$y = (y^{(1)}, y^{(2)}, \dots, y^{(m)})^T$$

コスト関数の最小化

コスト関数を最小化するようなパラメータ θ をどのように求めるか
入力が1変数の場合で考えてみる

コスト関数 $J(\theta) = J(\theta_0, \theta_1)$ はパラメータ θ_0, θ_1 について凸関数であり、 $J(\theta_0, \theta_1)$ が最小となるようなパラメータ θ_0^*, θ_1^* が見つけられればよい



線形回帰

勾配ベクトル

コスト関数 $J(\theta) = J(\theta_0, \theta_1)$ をパラメータ θ_0, θ_1 でそれぞれ偏微分した時のベクトル

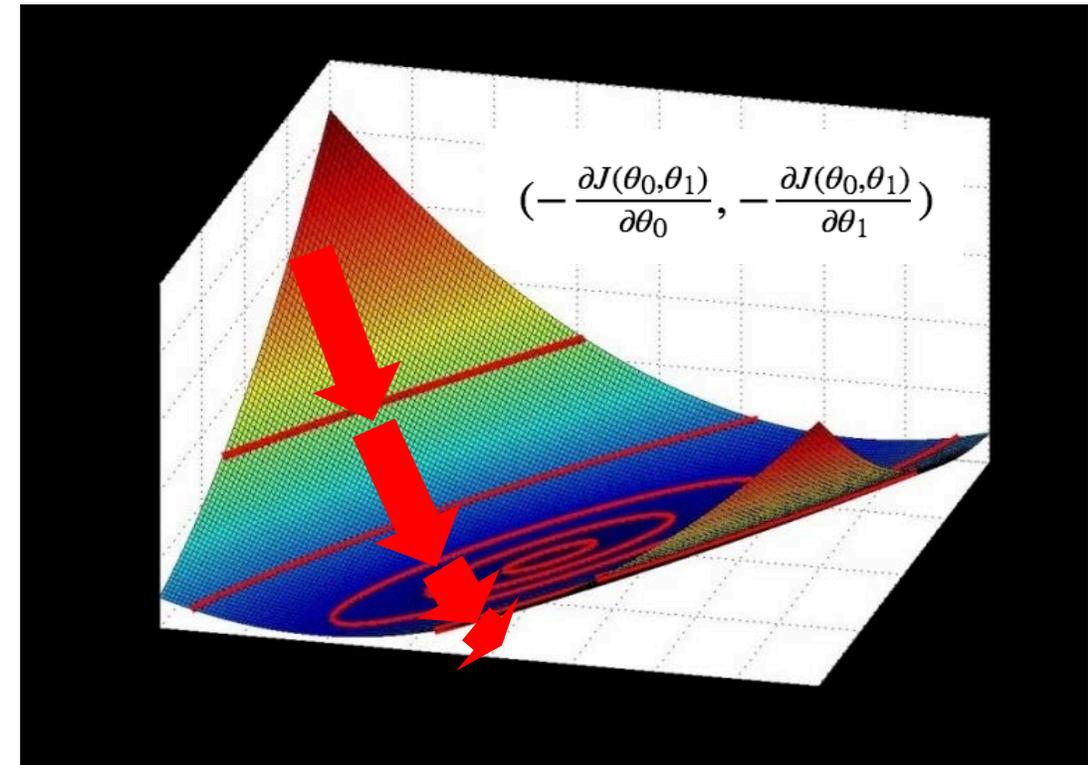
$$\nabla J(\theta_0, \theta_1) = \left(\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0}, \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} \right)$$

勾配ベクトルの方向：

- その点から移動した時に関数の値が最も大きく（逆方向は小さく）なる方向
 - 図の曲面の等高線の法線ベクトルとなっている

勾配ベクトルの大きさ：

- その点から勾配ベクトルの方向に移動した時に関数の値は $\|\nabla J(\theta_0, \theta_1)\|$ だけ増える（逆方向は減る）



線形回帰

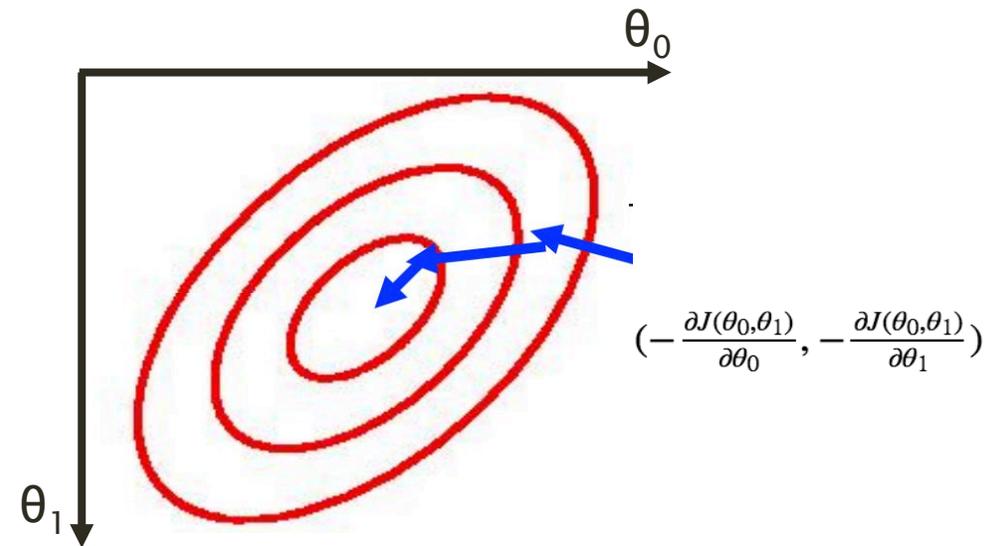
最急降下法

勾配ベクトルの定義より、 $-\nabla J(\theta_0, \theta_1)$ の方向に進むようにパラメータを更新すれば $J(\theta_0, \theta_1)$ を減少させることができる

この時、どれくらい勾配ベクトルの方向へ進むようかを定めるハイパーパラメータを α (学習率) として、パラメータ θ は任意の初期値から開始して $J(\theta_0, \theta_1)$ を減少させる方向に以下の様に更新できる

$$\theta_0 := \theta_0 - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} = \theta_0 - \alpha \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) / m$$

$$\theta_1 := \theta_1 - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} = \theta_1 - \alpha \sum_{i=1}^m ((h(x^{(i)}) - y^{(i)}) x^{(i)}) / m$$



補：勾配ベクトルの導出

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} = \frac{1}{2m} \sum_{i=1}^m \frac{\partial (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})$$

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} = \frac{1}{2m} \sum_{i=1}^m \frac{\partial (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) \frac{\partial (\theta_0 + \theta_1 x^{(i)} - y^{(i)})}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) x^{(i)}$$

最急降下法

入力の変数が多変数（ n 次元の特徴量ベクトル）の時も同様に以下の様に各パラメータ θ_j を更新してコスト関数 $J(\theta)$ が最小となるパラメータを求めればよい

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j} = \theta_j - \alpha \sum_{i=1}^m ((h(x^{(i)}) - y^{(i)}) x_j^{(i)}) / m$$

訓練データ $x^{(i)}$ の誤差の大きさに応じて j 次元のパラメータを調整

データセットの入力 X と正解 y を元にコスト関数 $J(\theta)$ を最小にするようなパラメータ θ を求める更新は以下のように表すことができる

$$\theta := \theta - \frac{\alpha}{m} X^T (X\theta - y)$$

$$X = \begin{pmatrix} x_0^{(1)} & x_1^{(1)} & \dots & x_n^{(1)} \\ \dots & \dots & \dots & \dots \\ x_0^{(m)} & x_1^{(m)} & \dots & x_n^{(m)} \end{pmatrix} \quad y = (y^{(1)}, y^{(2)}, \dots, y^{(m)})^T$$

$$\theta = (\theta_0, \theta_1, \dots, \theta_n)^T$$

最急降下法

バッチ勾配降下法

- データセット全体の誤差に基づく勾配によりパラメータを更新
 - パラメータ更新の計算が遅い
 - 大規模なデータセットではオンメモリで処理できない
 - 非凸なコスト関数では局所最適解に陥る可能性がある

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j} = \theta_j - \alpha \sum_{i=1}^m ((h(x^{(i)}) - y^{(i)})x_j^{(i)})/m$$

確率的勾配降下法

- ランダムに選択した各訓練データの誤差に基づく勾配によりパラメータを更新
 - パラメータ更新の計算が早い
 - 局所最適解から抜け出せる可能性がある
 - 最適解への収束が不安定
 - 学習率を徐々に小さくすることにより収束を安定させることができる

$$\theta_j := \theta_j - \alpha (h(x^{(i)}) - y^{(i)})x_j^{(i)}$$

最急降下法

訓練データセット： $\{(x^{(i)}, y^{(i)}); i=1, \dots, m\}$

1. 入力の各特徴量を標準化（必要であれば出力も標準化）
 - ・特徴量のスケールを揃えることで最急降下法の収束を早めることができる
2. 入力の各データ $x^{(i)}$ にバイアス項（ $x_0^{(i)}=1$ ）を追加
3. パラメータの初期値を与える
4. 訓練データセットを元にパラメータを更新
 - ・パラメータが収束するあるいは一定の繰り返し回数終了するまで
5. 学習されたパラメータを元にテストデータに対して予測を行いモデルを評価する
 - ・*学習モデルの選択と評価については次回説明

学習モデルの評価

訓練データセットで学習されたモデルのよさは、訓練データとは別に用意したテストデータに対して学習モデルを用いて予測を行い評価する

決定係数

- 学習モデルの予測値 \hat{y} と実際のデータ y との相関係数の2乗
 - 観測データの分散のうちどれぐらいを学習モデルが説明できているかを表す
 - 1に近いほどデータへの学習モデルのあてはまりがよい

$$R^2 = \frac{(\sum_{i=1}^m (\hat{y}^{(i)} - \bar{\hat{y}})(y^{(i)} - \bar{y}))^2}{(\sum_{i=1}^m (\hat{y}^{(i)} - \bar{\hat{y}})^2)(\sum_{i=1}^m (y^{(i)} - \bar{y})^2)} = \frac{\sum_{i=1}^m (\hat{y}^{(i)} - \bar{\hat{y}})^2}{\sum_{i=1}^m (y^{(i)} - \bar{y})^2}$$

ただし $\bar{\hat{y}} = \bar{y}$

テストデータに対する残差平方和や十分小さい・決定係数が十分大きいならよい

一方、訓練データに学習モデルが適合しすぎてしまい、テストデータのような未知のデータでは予測がうまくいかなくなることがある（過学習の問題）

正則化などにより過学習を防ぎ、なるべく未知のデータへもあてはまるようにしたい（汎化）

正規方程式

コスト関数を最小化するようなパラメータは解析的に求めることが可能

コスト関数の勾配が0になるようなパラメータを求める

$$J(\theta) = \frac{1}{2m}(X\theta - y)^T(X\theta - y) \quad \nabla J(\theta) = 0$$

以下の正規方程式が得られる

- 導出にはトレースを用いた行列の微分を使う

$$\theta = (X^T X)^{-1} X^T y$$

$X^T X$ の逆行列が計算できればコスト関数を最小化する最適パラメータを解析的に求められる

$X^T X$ が正則（フルランク）であること（階数 $\text{rank}(X^T X)$ が最大であること）

データ数 \ll 特徴量数の時は解が一意に定まらない（計算も大変 $O(n^3)$ ）

- 正則でない場合は $X^T X$ の対角成分に正の定数を加えて正則にする（過学習に対する正則化に関連）

多重共線性

入力データの特徴量間に強い相関があるような場合は注意が必要

正規方程式によりパラメータを求める場合の $X^T X$ の正則性に関連

- $X^T X$ が正則であるには、 X の各列（特徴量）が1次独立である必要がある

具体的な問題としては

- 相関がある特徴量それぞれにパラメータの重み与えようとしてパラメータの分散が大きくなる
- 特にパラメータの重みにより学習モデルを解釈する場合は問題になる

お互いに相関が強い特徴量があれば特徴量選択などによりいずれかの特徴量を採用するなどする

復習：正規分布

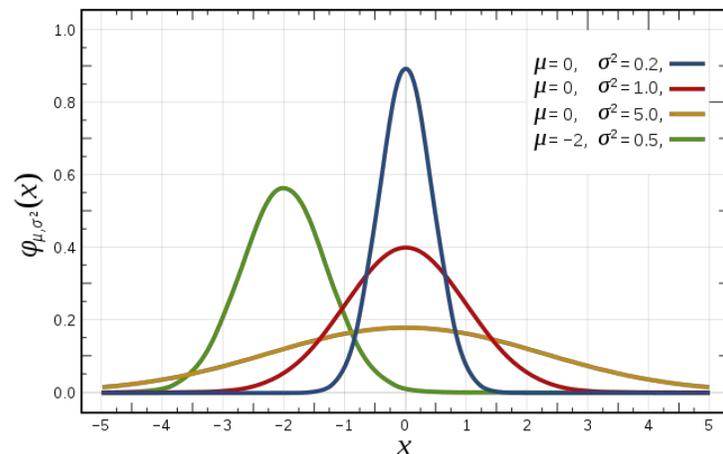
確率変数が正規分布に従うとき、その確率密度関数は

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \quad (x \in \mathbb{R})$$

正規分布に従う確率変数の平均（期待値） $E(X)$ 、分散 $V(X)$ は

- $E(X)=\mu$, $V(X)=\sigma^2$

正規分布は平均、分散の2つの母数で決まる



確率変数の平均（期待値）と分散

- 離散型
 - 平均
 - 確率変数のとりうる値とその確率の積

$$E(X) = \sum_i x_i P(X = x_i)$$

- 分散

$$V(X) = \sum_i (x_i - E(X))^2 P(X = x_i)$$

- 連続型

- 平均

$$E(X) = \int_{-\infty}^{\infty} x f(x) dx$$

- 分散

$$V(X) = \int_{-\infty}^{\infty} (x - E(X))^2 f(x) dx$$

復習：最尤法によるパラメータ推定

最尤法によるパラメータ推定

- 最尤原理
 - 現実の標本は確率最大のものが実現した
- 尤度関数
 - パラメータ θ の元で、観測したデータがどの程度起こりうるかを表す関数
 - 確率分布（確率密度関数）の積をパラメータ θ の関数とみなしたもの

$$L(\theta) = \prod_{i=1}^n f(x_i|\theta)$$

- 尤度関数を（パラメータがとる値の集合の空間で）最大にするものを観測したデータに対するパラメータの推定量とする
- 対数尤度
 - 尤度関数の対数をとって和の形にしたもの
 - 観測したデータについて対数尤度を最大にするパラメータの推定：最尤推定

$$\log L(\theta) = \sum_{i=1}^n \log f(x_i|\theta) \quad \frac{\partial \log L(\theta)}{\partial \theta} = 0$$

線形回帰

最尤推定による解釈

訓練データセットの出力 $y^{(i)}$ は平均 $h(x^{(i)})$ 、分散 σ^2 の正規分布に従う確率変数の実現値であるとする

- $y^{(i)}$ は $h(x^{(i)})$ に平均0、分散 σ^2 の正規分布に従うようなノイズ $\epsilon^{(i)}$ が合わさり観測されるとも考えられる
 - $y^{(i)} = h(x^{(i)}) + \epsilon^{(i)}$
- $\epsilon^{(i)}$ は独立同分布と仮定する

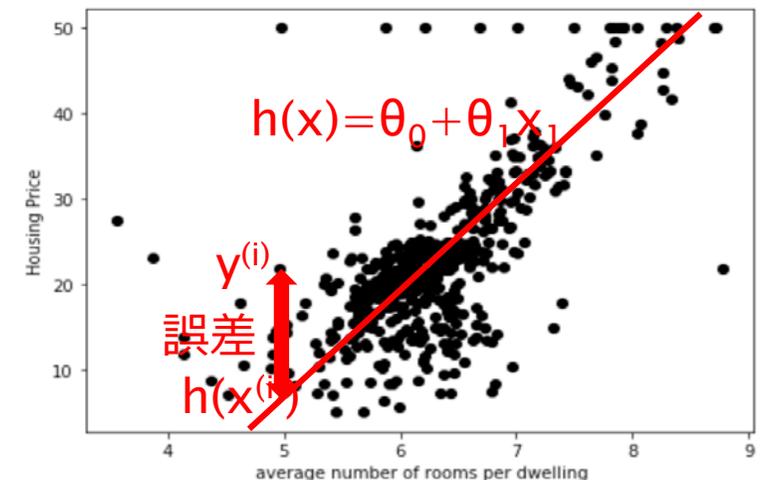
$\epsilon^{(i)}$ の確率密度関数（データ $x^{(i)}$ が与えられた時の $y^{(i)}$ の分布（パラメータは θ ））は

$$f(\epsilon^{(i)}) = f(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - h(x^{(i)}))^2}{2\sigma^2}\right)$$

最尤推定によりパラメータをする

対数尤度は

$$\log L(\theta) = \log \prod_{i=1}^m f(y^{(i)} | x^{(i)}; \theta) = \log \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - h(x^{(i)}))^2}{2\sigma^2}\right)$$



最尤推定による解釈

対数尤度 (つづき)

$$\begin{aligned} \log L(\theta) &= \log \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - h(x^{(i)}))^2}{2\sigma^2}\right) \\ &= \sum_{i=1}^m \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - h(x^{(i)}))^2}{2\sigma^2}\right) \\ &= m \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} \sum_{i=1}^m (y^{(i)} - h(x^{(i)}))^2 \end{aligned}$$

第1項は定数、対数尤度を最大にするには第2項を最小化すればよい

コスト関数に用いた残差平方和を最小化することに等しい $\sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$

最小二乗法による線形回帰はパラメータの最尤推定をしているとも考えられる