

# 統計的モデリング基礎⑥ ～正則化と事後確率最大化～

鹿島久嗣  
(情報学科 計算機科学コース)

# 計算グラフと自動微分

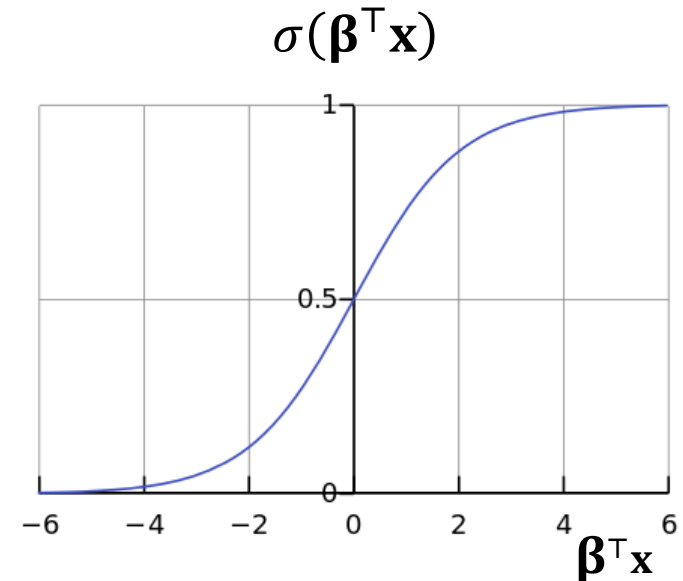


# ロジスティック回帰： ダミー変数を従属変数とするモデル

- 以前、重回帰モデルでダミー変数を従属変数とすると、厳密には少しおかしいという話だった → もっときちんと扱いたい
  - 重回帰モデル  $y = \boldsymbol{\beta}^T \mathbf{x}$  の従属変数の値域は実数全体
- 従属変数の値域が  $\{-1, +1\}$  もしくは  $(0, 1)$  ( $Y = +1$  となる確率) となるようにしたい
- ロジスティック回帰モデル：

$$P(Y = 1 | \mathbf{x}, \boldsymbol{\beta}) = \frac{1}{1 + \exp(-\boldsymbol{\beta}^T \mathbf{x})} = \sigma(\boldsymbol{\beta}^T \mathbf{x})$$

- $\sigma$ ：ロジスティック関数 ( $\sigma: \mathbb{R} \rightarrow (0, 1)$ )



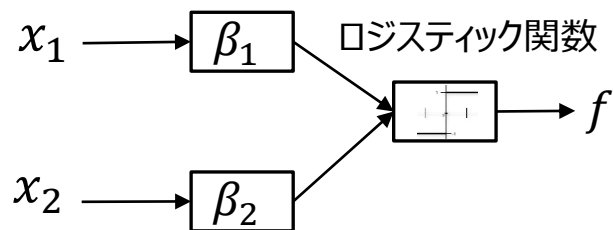


# ニューラルネットワーク：

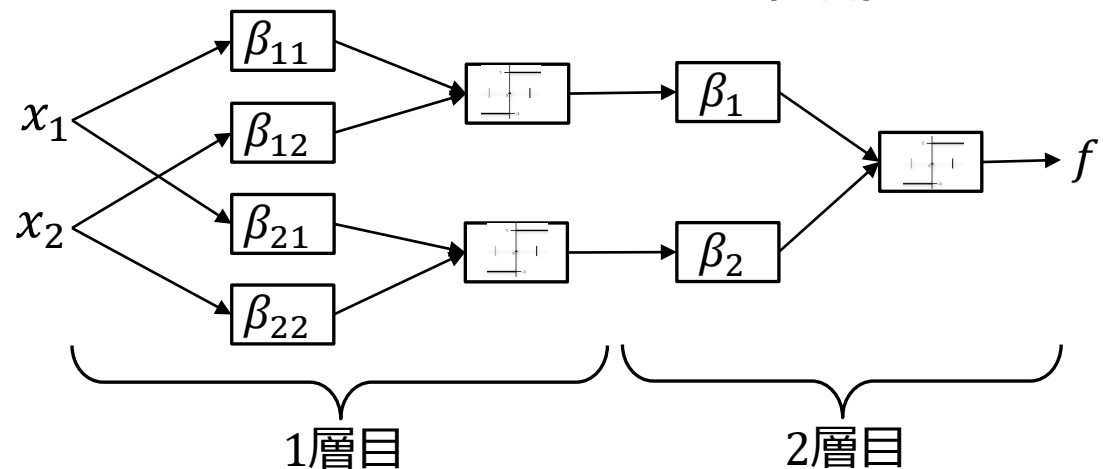
(ざっくりいえば) ロジスティック回帰モデルを連結したものの

- ニューラルネットワークはロジスティック回帰モデルを連結したものの
  - 複数のロジスティック回帰モデルの出力が、別のロジスティック回帰モデルの入力になる
  - ロジスティック関数（非線形）によりモデルに非線形性を導入
  - 両者ともに、 $y = +1$ である確率 $f(\mathbf{x}; \boldsymbol{\beta})$ を出力するモデル

ロジスティック回帰モデル



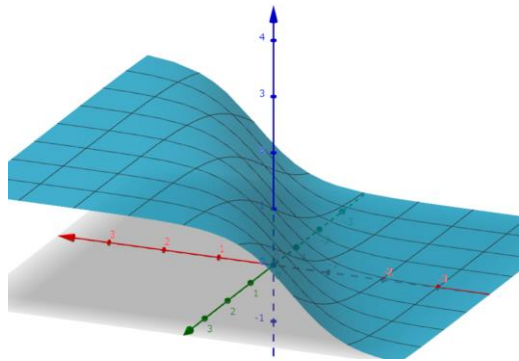
ニューラルネットワーク（2層）



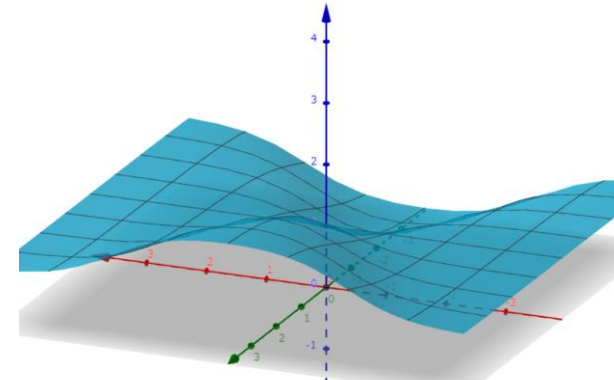
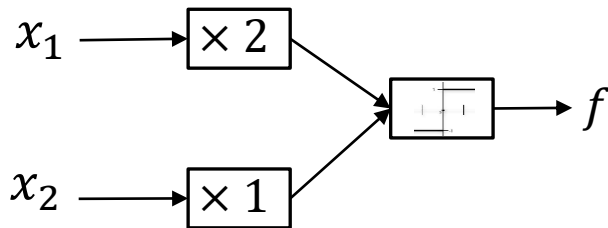


# ニューラルネットワークの非線形性の例： ロジスティック回帰を2層積むと非線形分類が可能

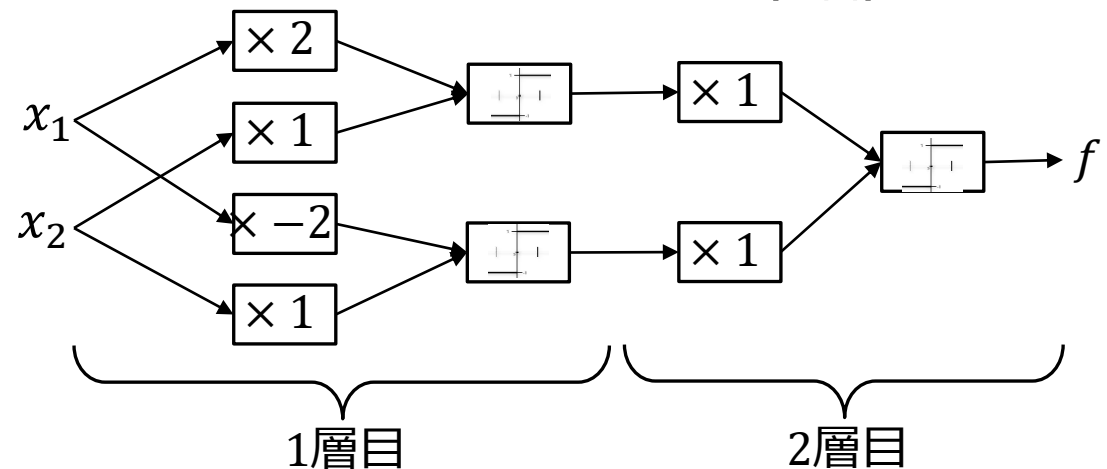
- ロジスティック回帰は1層では線形判別しかできない (AND/OR)
- 2層以上積むことで非線形の表現力を獲得 (XOR)



ロジスティック回帰モデル



ニューラルネットワーク (2層)





# ニューラルネットワークのパラメータ推定： 最急降下法を適用するために勾配の計算が必要

- 対数尤度関数 $L(\boldsymbol{\beta})$ を最大化するパラメータ $\boldsymbol{\beta}$ を求める：

$$L(\boldsymbol{\beta}) = - \sum_{i=1}^n \left( \delta(y^{(i)} = 1) \log f(x^{(i)}) + \delta(y^{(i)} = -1) \log (1 - f(x^{(i)})) \right)$$

- $f(x^{(i)})$ は $x^{(i)}$ に対するニューラルネットの出力 ( $y^{(i)} = 1$ である確率)
- 勾配 $\nabla L(\boldsymbol{\beta}) = \frac{\partial L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}}$ が計算できれば最急降下法を適用できる：  
$$\boldsymbol{\beta}^{\text{NEW}} \leftarrow \boldsymbol{\beta} + \eta \nabla L(\boldsymbol{\beta})$$
  - 実際は確率的最適化やミニバッチを用いることも多い

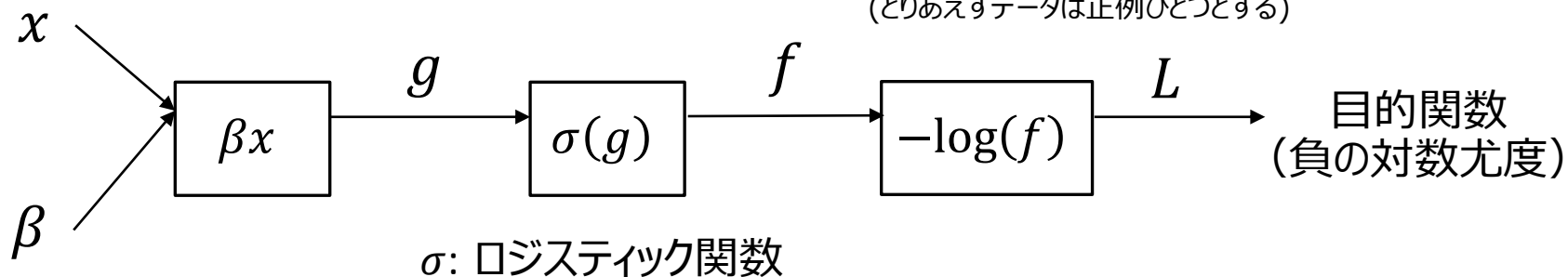


# 計算グラフ：

ニューラルネットの入力から出力までの計算を図示

- 計算グラフ：関数の入出力の間を、単純な計算ユニットをつないで表したもの
- 計算グラフをたどりながら、入力に順番に単純な操作（重み付き和やロジスティック関数の適用など）を適用していくと、出力が得られる

- **ロジスティック回帰の計算グラフ**：  
ロジスティック回帰の出力： $f = \sigma(\beta x)$   
最適化問題の目的関数： $L = -\log f$   
(とりあえずデータは正例ひとつとする)



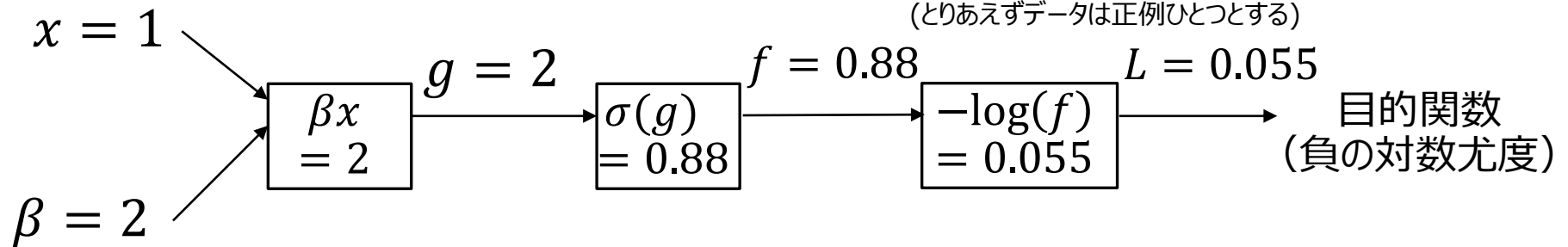


# 計算グラフ：

ニューラルネットの入力から出力までの計算を図示

- 計算グラフ：関数の入出力の間の関係を、単純な計算ユニットをつないで表したもの
- 計算グラフをたどりながら、入力に順番に単純な操作（重み付き和やロジスティック関数の適用など）を適用していくと、出力が得られる

- **ロジスティック回帰の計算グラフ**：ロジスティック回帰の出力： $f = \sigma(\beta x)$   
最適化問題の目的関数： $L = -\log f$   
(とりあえずデータは正例ひとつとする)



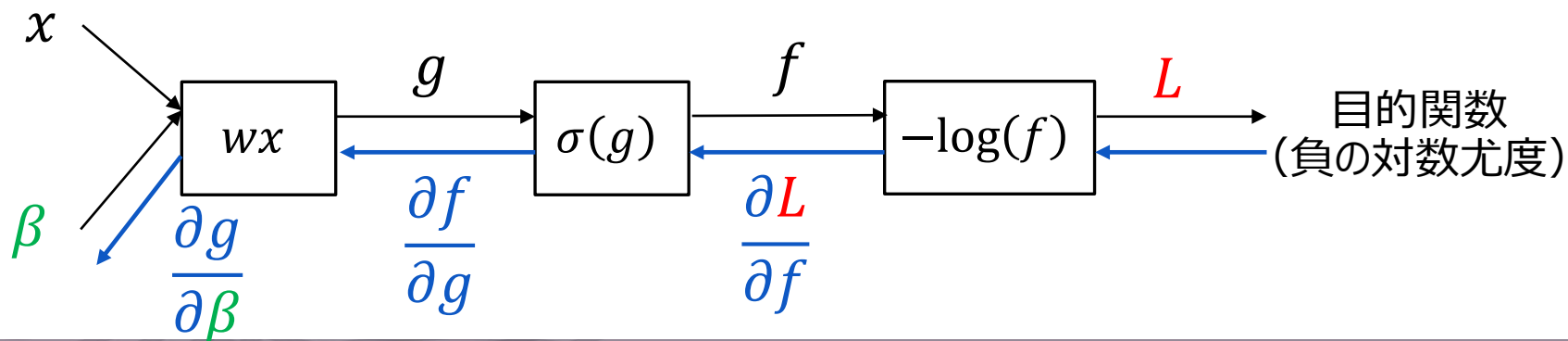
$\sigma$ : シグモイド関数

# 計算グラフ上での自動微分：

計算グラフを出力から逆向きにたどることで勾配計算

- 勾配計算： $\partial L / \partial \beta$ を求めたい
- 計算グラフ上で $L$ と $\beta$ は遠い
- 計算グラフを逆向きにたどりながら、微分を計算する

ー ロジスティック回帰の場合：
$$\frac{\partial L}{\partial \beta} = \frac{\partial g}{\partial \beta} \cdot \frac{\partial f}{\partial g} \cdot \frac{\partial L}{\partial f}$$

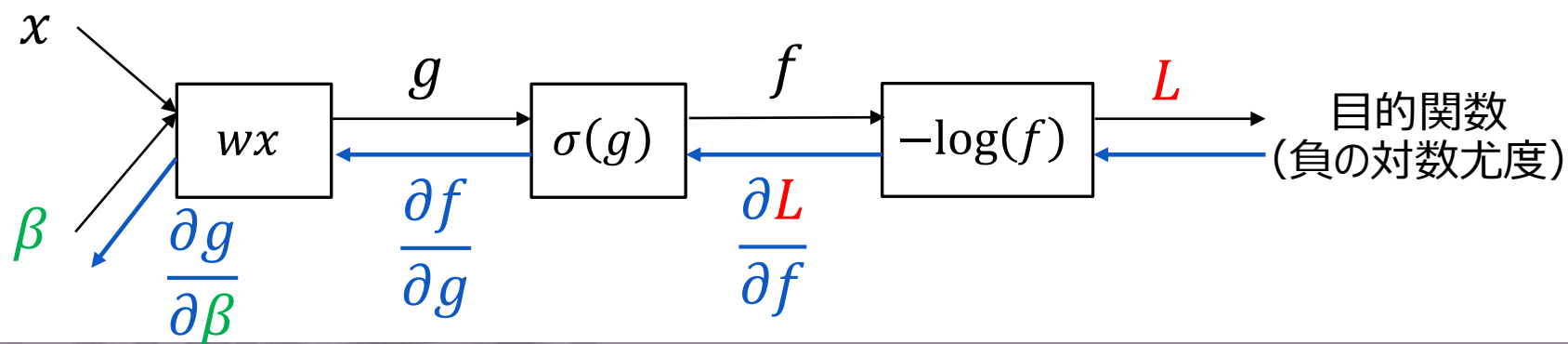


# 自動微分のポイント： 微分可能な計算ユニット

- 自動微分：計算グラフを逆向きにたどりながら（微分の連鎖律によって）微分を計算する
- 各ユニットは、入力について微分可能である必要がある

- $\frac{\partial f}{\partial g} = \frac{\partial \sigma(g)}{\partial g} = \sigma(g)(1 - \sigma(g))$ （ロジスティック関数の微分）

- $\frac{\partial g}{\partial \beta} = \frac{\partial \beta x}{\partial \beta} = x$



## ニューラルネットワーク推定のポイント：

微分可能なユニットを組みあわせて自動微分にまかせる

- ニューラルネットワーク推定法の汎用性
  1. ネットワークを「計算グラフ」で記述する
    - 各ユニットはパラメータや入力について微分可能とする
  2. 誤差逆伝播で自動的に勾配が計算できる  
(自動微分)

# 正則化

# 重回帰モデルの復習： 最小二乗法による定式化

■ 重回帰モデル： $y = \boldsymbol{\beta}^\top \mathbf{x}$

• パラメータ： $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_m, \alpha)^\top$

• 独立変数： $\mathbf{x} = (x_1, x_2, \dots, x_m, 1)^\top$

} 最後の次元は  
切片部分に相当

■ データ：

• 計画行列： $\mathbf{X} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)})^\top$

• 従属変数： $\mathbf{y} = (y^{(1)}, y^{(2)}, \dots, y^{(n)})^\top$

■ 目的関数： $\ell(\boldsymbol{\beta}) = \sum_{i=1}^n (y^{(i)} - \boldsymbol{\beta}^\top \mathbf{x}^{(i)})^2 = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$

# 重回帰モデルの解： 解析解が得られる

- 目的関数： $L(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$
- 解： $\boldsymbol{\beta}^* = \operatorname{argmin}_{\boldsymbol{\beta}} L(\boldsymbol{\beta}) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$
- ただし、解が存在するためには $\mathbf{X}^\top \mathbf{X}$ が正則である必要
  - モデルの次元数 $m$ よりもデータ数 $n$ が大きい場合はおおむね成立
- 正則化：正則でない場合には $\mathbf{X}^\top \mathbf{X}$ の対角成分に正の定数 $\lambda > 0$ を加えて正則にする
  - 新たな解： $\boldsymbol{\beta}^* = \operatorname{argmin}_{\boldsymbol{\beta}} L(\boldsymbol{\beta}) = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$
  - 目的関数に戻ると： $L(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_2^2$

パラメータのノルムに関する  
ペナルティ項

## データへの過適合：

データへの過剰な適合は将来のデータへの予測力を損なう

- 先ほどは、正則化を計算を安定させるために導入した
- 例えば、データ数 $n$ が次元数 $m$ より小さいとき、重回帰の解は一意に定まらない
  - 任意の数の解が存在し、どれが良いのかわからない
- データへの過適合：
  - 汎化：予測を目的とする場合、我々の真の目的は将来のデータへの予測力が高いモデルを得ること
  - 手持ちのデータへのモデルの過剰な適合は、将来の予測力を損なう可能性がある



# オッカムの剃刀：

できるだけ“単純な”なモデルを採用せよ

---

- データに同程度適合している無数のモデルのうち、どれが最も“良い”モデルだろうか
- オッカムの剃刀：単純なモデルを採用せよ
  - 「ある事柄を説明するためには、必要以上に多くを仮定するべきでない」
- 単純さを何で測るか？
  - 例えば、モデルに含まれる独立変数の数
  - 自由度調整済決定係数、AIC、BICなどの情報量基準

## 0-ノルムを用いた正則化：

パラメータ中の非零成分の数を減らす

- モデルに含まれる独立変数の数 =  $\beta$  中の非零成分の数 =  $\beta$  の 0-ノルム

モデルに含まれる  
独立変数の数

- 0-ノルム制約を入れた回帰問題：

$$\beta^* = \operatorname{argmin}_{\beta} \|y - X\beta\|_2^2 \text{ s.t. } \|\beta\|_0 \leq \eta$$

- あるいは 0-ノルムをペナルティ項として導入：

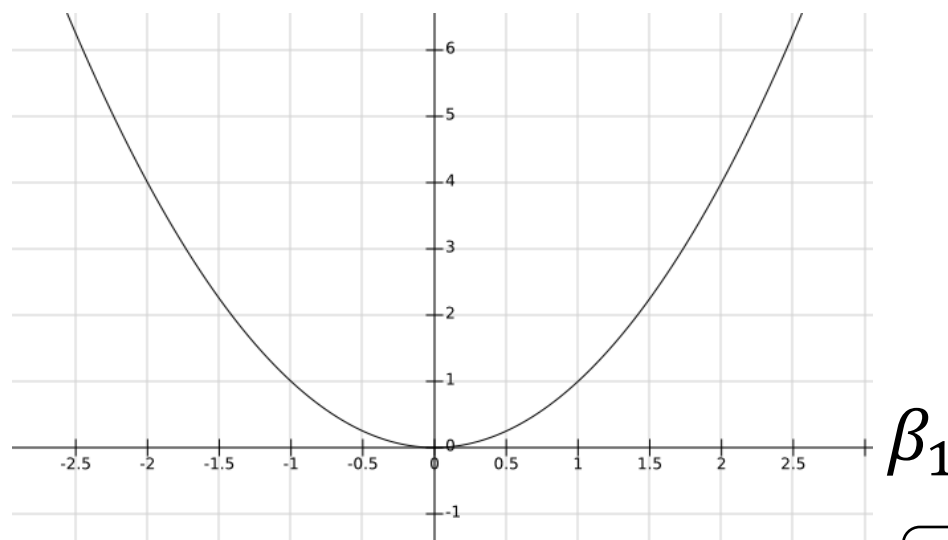
$$\beta^* = \operatorname{argmin}_{\beta} \|y - X\beta\|_2^2 + \lambda \|\beta\|_0$$

- $\eta$  と  $\lambda$  は一対一対応している

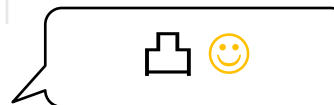
- ただし、この問題は凸最適化問題でない

# 0-ノルムの代替： 2-ノルム正則化は凸最適化になる

- 0-ノルム  $\|\beta\|_0$  の代わりに 2-ノルム  $\|\beta\|_2^2$  を用いる



- リッジ回帰： $L(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_2^2$



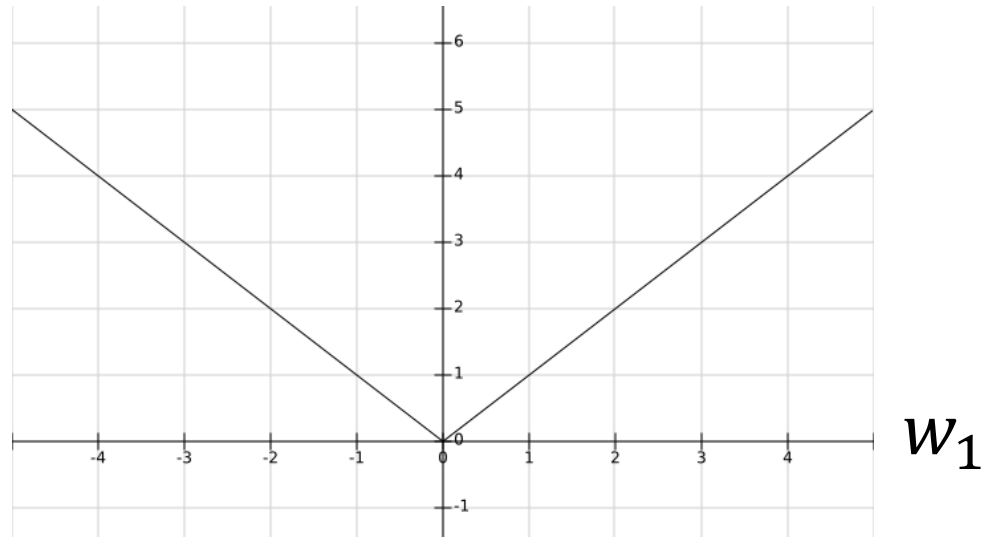
- 0-ノルム正則化の緩和版として捉えることができる：

$$L(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_0$$



# 0-ノルムの代替： 1-ノルム正則化も凸最適化になる

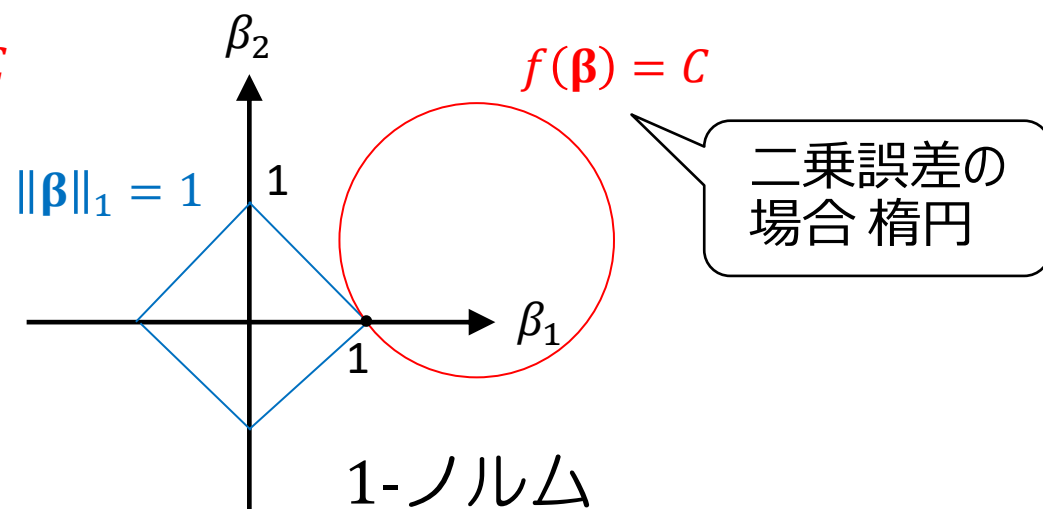
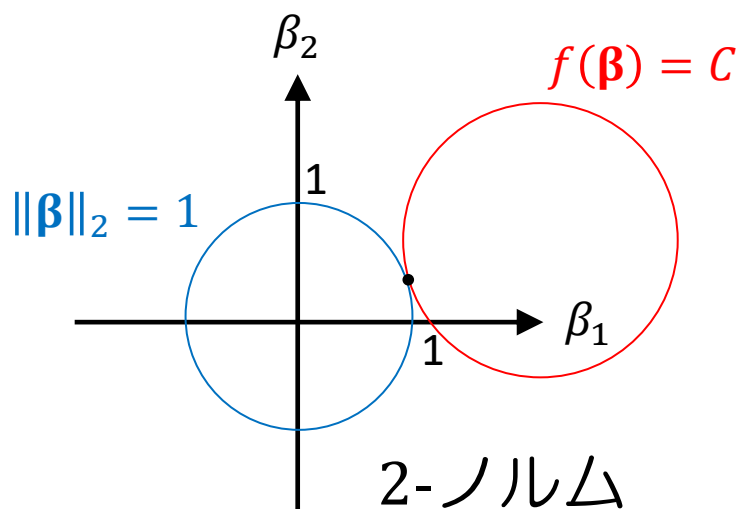
- さらに、1-ノルム  $\|\boldsymbol{\beta}\|_1 = |\beta_1| + |\beta_2| + \dots + |\beta_D|$  も利用可能



- ラツソ：  $L(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_1$ 
  - 凸最適化だが、解析解はもたない
- 1-ノルムを用いると疎な解になる ( $\boldsymbol{\beta}^*$ の多くの要素が0になる)

# 1-ノルム正則化の利点： 疎な解をもつ

- 1-ノルム正則化は1-ノルム制約で書き直せる：  
$$\operatorname{argmin}_{\boldsymbol{\beta}} f(\boldsymbol{\beta}) + \gamma \|\boldsymbol{\beta}\|_1 \Leftrightarrow \operatorname{argmin}_{\boldsymbol{\beta}} f(\boldsymbol{\beta}) \text{ s.t. } \|\boldsymbol{\beta}\|_1 \leq \lambda$$
- 1-ノルム正則化は疎な解をもつ
  - 2-ノルム制約の等高線（円形）と1-ノルム制約の等高線（菱形）の比較



二乗誤差の  
場合 楕円

# 事後確率最大化推定

# 回帰のベイズ統計的解釈： 事後確率最大化推定

---

- 線形回帰モデルの推定は最尤推定として解釈できた
- 正則化のもとでの回帰モデルの推定はベイズ統計の枠組みで解釈できる
  - 事前分布・事後分布の導入
  - 事後確率最大化（MAP）推定
  - リッジ回帰 = MAP推定

# 線形回帰モデルの最尤推定：

線形回帰の確率モデルの最尤推定 = 最小二乗法

- 線形回帰モデルに対応する確率モデル

- 確率密度関数：
$$f(y^{(i)} | \mathbf{x}^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \boldsymbol{\beta}^\top \mathbf{x}^{(i)})^2}{2\sigma^2}\right)$$

- 対数尤度：
$$L(\boldsymbol{\beta}) = \sum_{i=1}^n \log f(y^{(i)} | \mathbf{x}^{(i)})$$
$$= -\frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \boldsymbol{\beta}^\top \mathbf{x}^{(i)})^2 + \text{const.}$$

- 対数尤度を $\boldsymbol{\beta}$ について最大化すること（最尤推定）  
= 二乗誤差を $\boldsymbol{\beta}$ について最小化すること（最小二乗法）



# ベイズ的統計モデリングの考え方： 尤度の代わりに事後分布を考える

- 最尤推定（MLE）では尤度を最大化するパラメータ $\beta$ を求めた：

$$P(\mathbf{y} \mid \mathbf{X}, \beta) = \prod_{i=1}^n f(y^{(i)} \mid \mathbf{x}^{(i)}, \beta)$$

- これは、パラメータが与えられたもとでデータが再現される条件付確率： $P(\text{データ} \mid \text{パラメータ}) = P(\mathbf{y} \mid \mathbf{X}, \beta)$

※ 回帰の場合は $\mathbf{X}$ が定数として与えられ、 $\mathbf{y}$ のみが確率変数と思っている

- ベイズ的なモデリングの考え方では、事後分布を考える：

$$P(\text{パラメータ} \mid \text{データ}) = P(\beta \mid \mathbf{X}, \mathbf{y})$$

- 事後分布はパラメータを確率変数と考える

# 事後分布：

事後分布  $\propto$  尤度  $\times$  事前分布

## ■ 事後分布：

$$P(\text{パラメータ} | \text{データ}) = \frac{P(\text{データ} | \text{パラメータ})P(\text{パラメータ})}{P(\text{データ})}$$

ベイズの定理

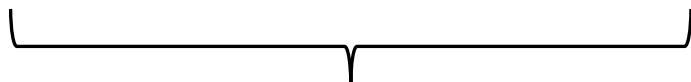
パラメータに依存しない  
(定数扱い)

- $P(\text{データ}) = \sum_{\text{パラメータ}} P(\text{データ} | \text{パラメータ})P(\text{パラメータ})$

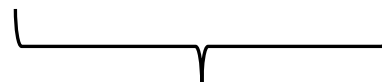
## ■ 対数事後分布：

$$\log P(\text{パラメータ} | \text{データ})$$

$$= \log P(\text{データ} | \text{パラメータ}) + \log P(\text{パラメータ}) + \text{const.}$$



尤度



事前分布

# 事後確率最大化 (MAP) 推定： 事後確率を最大化するパラメータを採用

## ■ 対数事後分布：

$$\log P(\text{パラメータ} \mid \text{データ})$$

$$= \log P(\text{データ} \mid \text{パラメータ}) + \log P(\text{パラメータ}) + \text{const.}$$

## ■ 事後確率最大化 (Maximum a posteriori; MAP) 推定

### ● 事後確率を最大化するパラメータを採用する：

$$\text{パラメータ}^* = \operatorname{argmax}_{\text{パラメータ}} \log P(\text{パラメータ} \mid \text{データ})$$

### ● 最尤推定では $\log P(\text{データ} \mid \text{パラメータ})$ の項のみ考える

### ● 追加の項として事前分布の項： $\log P(\text{パラメータ})$

# 事後確率最大化としてのリッジ回帰： 正規分布を事前分布とした事後確率最大化

- 対数事後分布：

$$\log P(\text{パラメータ} \mid \text{データ})$$

$$= \log P(\text{データ} \mid \text{パラメータ}) + \log P(\text{パラメータ}) + \text{const.}$$

- リッジ回帰：

対数尤度

対数事前分布

$$\boldsymbol{\beta}^* = \operatorname{argmin}_{\boldsymbol{\beta}} \frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \boldsymbol{\beta}^\top \mathbf{x}^{(i)})^2 + \frac{1}{2\sigma'^2} \|\boldsymbol{\beta}\|_2^2$$

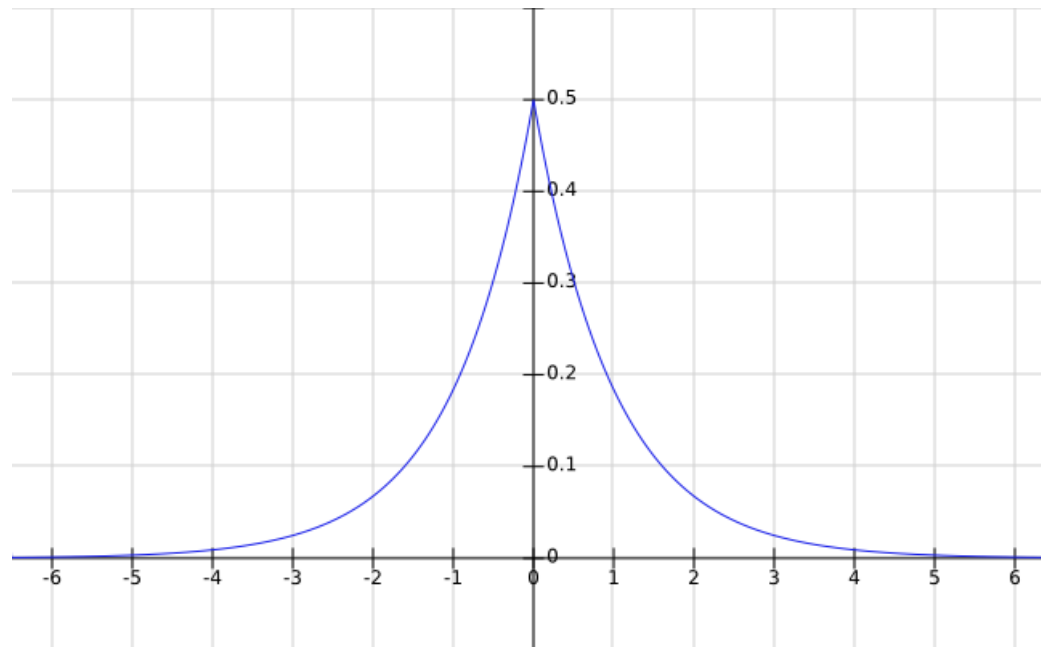
- ◆ 対数尤度： $\sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma'} \exp\left(-\frac{(y^{(i)} - \boldsymbol{\beta}^\top \mathbf{x}^{(i)})^2}{2\sigma'^2}\right)$

- ◆ 事前分布： $P(\boldsymbol{\beta}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\boldsymbol{\beta}^\top \boldsymbol{\beta}}{2\sigma^2}\right)$

# 事後確率最大化としてのラッソ： ラプラス分布を事前分布として利用

- 事前分布を正規分布にすると2-ノルム正則化
- ラプラス分布：1-ノルム正則化に対応する事前分布

$$P(\boldsymbol{\beta}) = \frac{1}{2\phi} \exp\left(-\frac{|\boldsymbol{\beta}|}{\phi}\right)$$



# ベイズ予測：

## 推定のばらつきを考慮した予測

- MAP推定では事後分布が最大となるパラメータを点推定する

$$\hat{\boldsymbol{\beta}} = \operatorname{argmax}_{\boldsymbol{\beta}} P(\boldsymbol{\beta} \mid \mathbf{X}, \mathbf{y})$$

- ベイズ予測では事後分布そのものを用いて予測する

$$P(y \mid \mathbf{x}) = \int_{\boldsymbol{\beta}} P(y \mid \mathbf{x}, \boldsymbol{\beta}) P(\boldsymbol{\beta} \mid \mathbf{X}, \mathbf{y}) d\boldsymbol{\beta}$$

- あらゆるパラメータにおけるモデルの予測を事後確率で重みづけて予測する
- 最適化問題を解いてパラメータを点推定するのではなく「全部」使う