

京都大学集中講義  
機械学習と深層学習の  
数理と応用  
(1)

鈴木大慈  
東京大学大学院情報理工学系研究科数理情報学専攻  
理研AIP  
Japan Digital Design

2018年12月17日

# 自己紹介

## ・現職：

- ・東京大学大学院情報理工学系研究科  
数理情報学専攻 准教授
- ・理研AIP「深層学習理論チーム」チームリーダー

## ・専門：

- ・機械学習の数理
  - ・汎化誤差理論
  - ・確率的最適化
- ・数理統計学
  - ・高次元統計
  - ・ノンパラメトリック法



# 本講義の目的

- ・対象：機械学習初学者
- ・機械学習の数理、特に深層学習の数理を解説
- ・Pythonによる簡単な実装も体験してもらう

## 評価

- ・出席とレポート
- ・レポート内容：
  - ・数学的問題に回答
  - ・Pythonによる深層学習の実装 (Google colab)

# 講義の予定

スライドを用いた講義

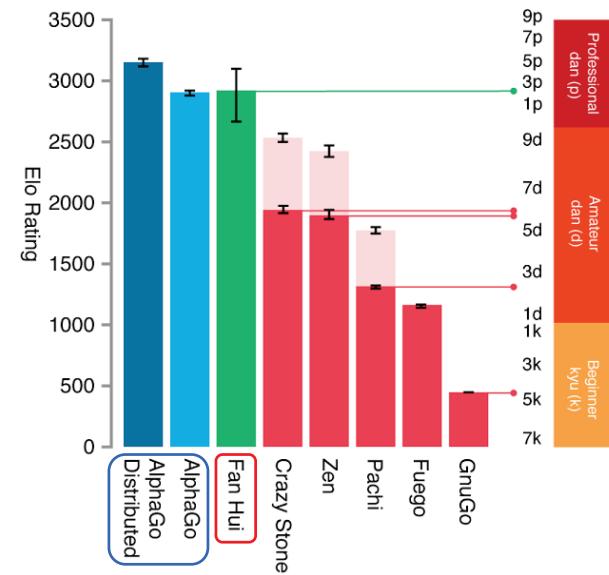
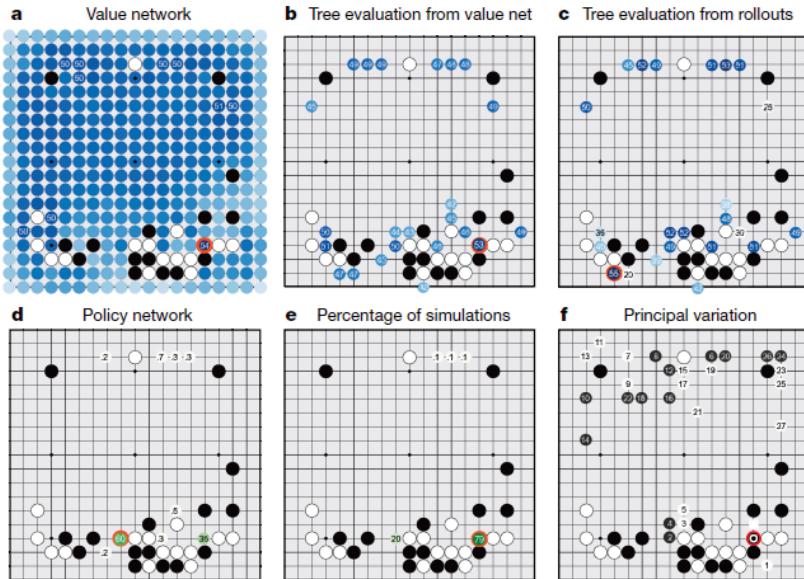
- 1回目：機械学習全般の概論
- 2回目：深層学習の概論

黒板を用いた講義

- 3回目：ニューラルネットワークの近似理論
  - 万能近似能力, Sobolev空間における近似精度
- 4回目：NNの汎化誤差解析
  - Rademacher複雑度
- 5回目：速い汎化誤差の収束レート, その他
  - Talagrandの不等式, ミニマックスレート最適性

# AlphaGo/Alpha Zero

深層学習+強化学習+モンテカルロ木探索+自己対戦



2015年10月5～9日 欧州覇者Fan Hui (2段) に5-0で勝利

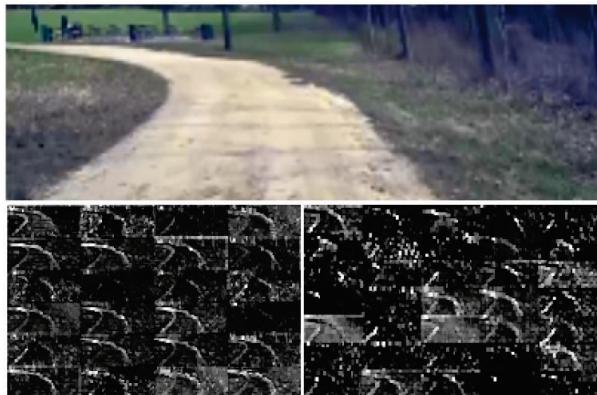
2016年3月9～15日 イ・セドルと対局し4-1で勝利

2017年5月23～27日 柯潔と対局し3-0で勝利

2017年12月 AlphaZeroが4時間の自己対戦のみでelmo, AlphaGo Zeroを上回る。

[Silver et al. (Google Deep Mind): Mastering the game of Go with deep neural networks and tree search, Nature, 529, 484–489, 2016]

# 自動運転



5月22日号 連動特集 エヌビディア、GS、トヨタ、日立  
ARTIFICIAL INTELLIGENCE  
世界制覇の攻防  
写真: Just\_Super/Getty Images

## 詳報：トヨタが頼った謎のAI半導体メーカー

産業秩序が激変、自動車を「操る」のは誰だ



島津 翔

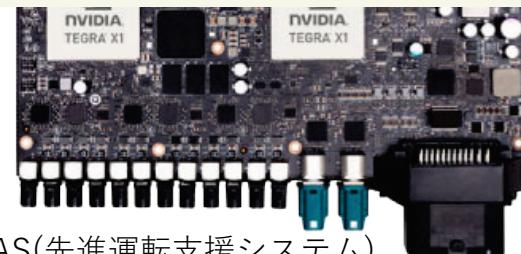
バックナンバー

2017年5月22日 (月)



AI（人工知能）による産業構造の激変が始まった。

売り上げ規模など従来の序列は全く関係ない。対応できない既存勢力は没落する。強固なピラミッドを持つ自動車産業で安泰ではない。AIによる自動運転の実用化が、激変の号砲



ADAS(先進運転支援システム)

[End to End Learning for Self-Driving Cars, Nvidia 2016]

# 身边にあふれる機械学習

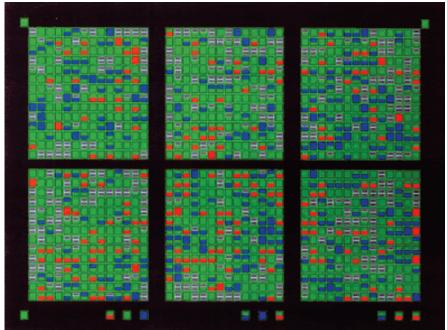
検索エンジン



推薦システム

The screenshot shows the Amazon Japan homepage with a search bar and navigation menu. A prominent red banner at the top says 'マイストア、おすすめ商品' (My Store, Recommended Products). Below it, a section titled '本日のみ' (Today Only) features a link to 'おすすめ商品を見る' (View Recommended Products). To the right, a large box displays recommended products based on the user's purchases, with a specific item highlighted: 'トワイニング クオリティアールグレイ 100g' (Twinings Quality English Breakfast 100g) for ¥864. The page also includes a rating section with stars and a purchase button.

遺伝子データ解析



音声認識



機械学習プラットフォーム

Google Cloud Platform

Microsoft Azure

amazon  
web services™

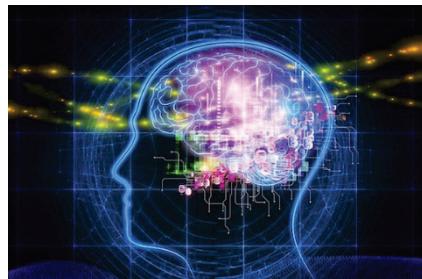
# 人工知能と機械学習

## 人工知能

# 本日の様相 「人工知能」 ≐ 「機械学習」

自分で問題設定ができ、  
その解決もできる。

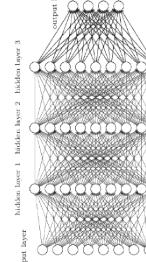
“汎用人工知能”



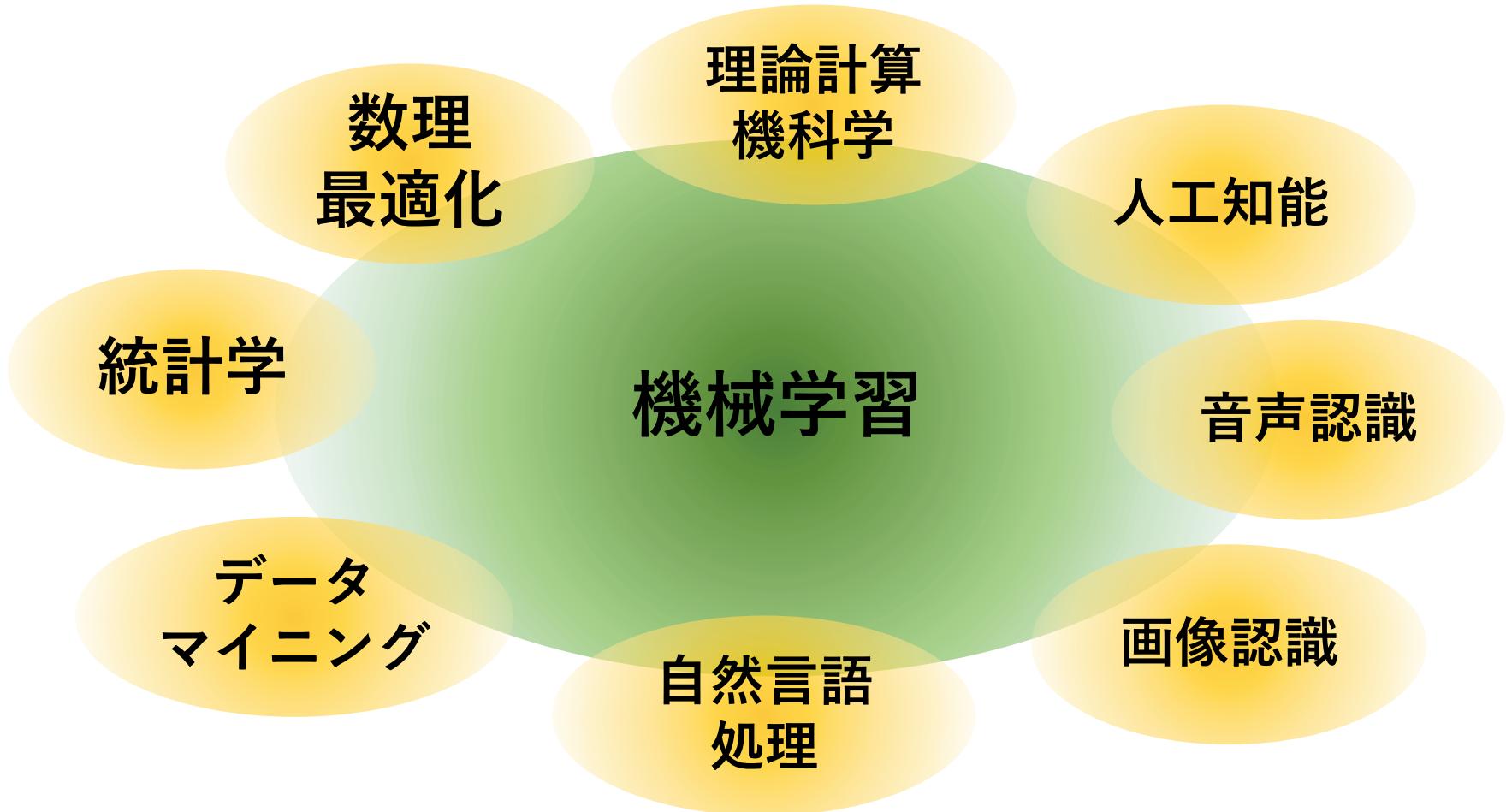
機械学習  
統計的アプローチ

SVM  
トピックモデル  
スペース学習  
テンソル学習  
...

深層学習



# 機械学習の立ち位置



さまざまな分野の複合領域

# 機械学習コミュニティの実体

## 機械学習主要国際会議

NeurIPS (Neural Information Processing Systems)

ICML (International Conference of Machine Learning)

COLT (Conference of Learning Theory)

ICLR (International Conference on Learning Representations)

AISTATS, UAI, ECML, …

## 関連国際会議

- データマイニング  
KDD, ICDM, WWW, WISDM, SIGIR, SDM
- 人工知能  
IJCAI, AAAI
- コンピュータビジョン  
CVPR, ICCV, ECCV
- 自然言語処理  
ACL, NAACL, EMNLP, COLING

➤ 全て査読あり：ダブルブラインド，採択率20～25%

- NIPS2016 (568/2500, 22.7%), ICML2016 (322/1327, 24.3%)



NIPS2015@Montreal



ICML2016@NYC

# NeurIPS2018

- ・投稿数：4856件
- ・採択数：1011件
- ・査読プロセス：double blind

(Senior) Area Chair: 約280人, 査読者：約2500人  
 ✓ 1査読者あたり6本を査読, そのうち1-2本がaccept

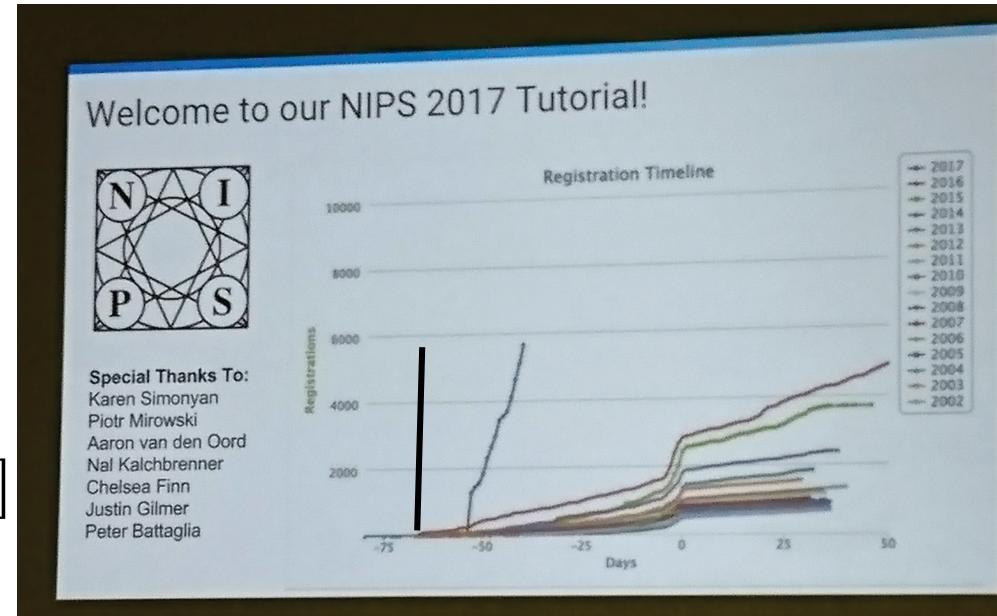


通常レジストレーション  
 (2000席)

11分38秒で完売

(論文採録者, トップ査読者は別)

[去年は8000席が2週間で完売]



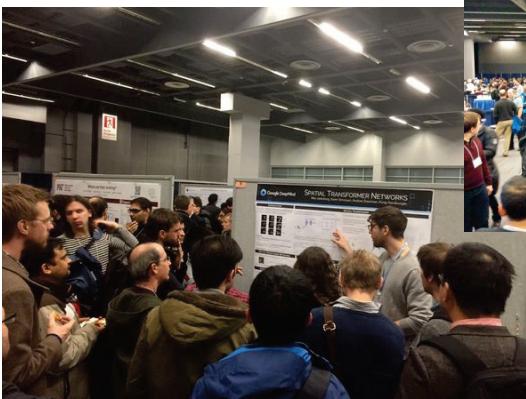
# 会場の様子



2015年の様子



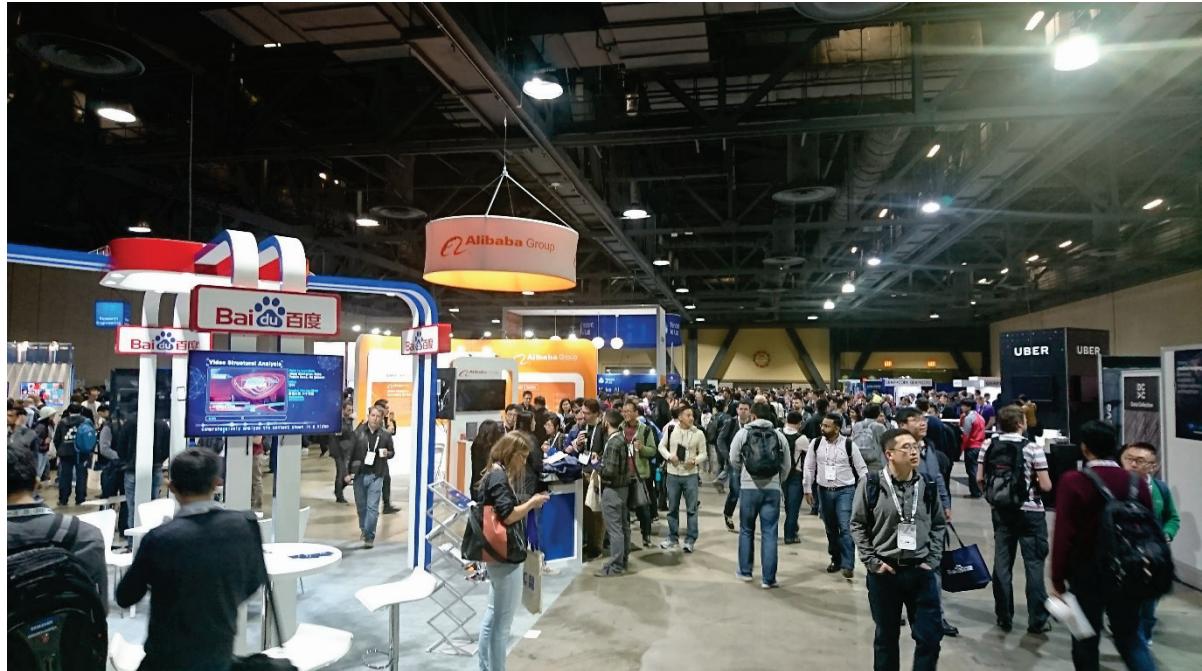
Deep learning tutorial



ポスター発表



# 企業ブースの様子 (NIPS2017)

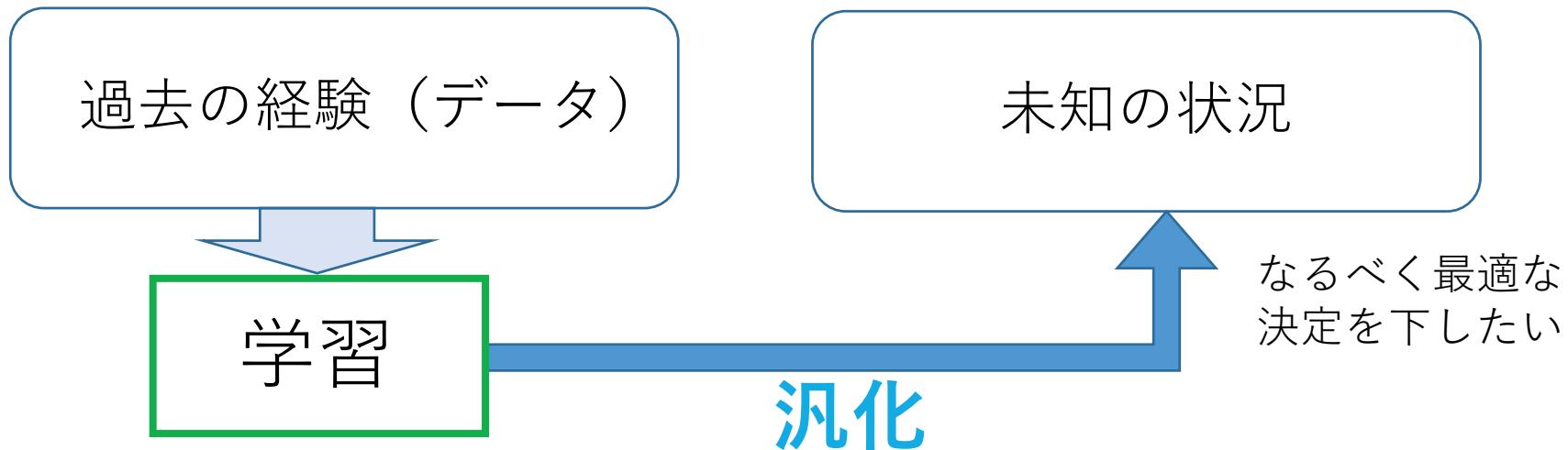


# 機械学習の目的

- 人間と同様の知的情報処理を計算機で実現するための技術・手法



Arthur Samuel 「Field of study that gives computers the ability to learn without being explicitly programmed」 (1959)



# 予測と推測

機械学習の活用法

## 予測

(より機械学習的)



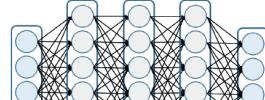
→ 船



→ “Hello”

- Outcomeを正しく当てる.
- 解釈よりも予測精度を重視.

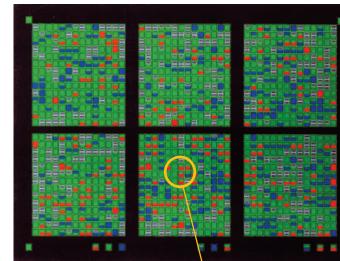
例：深層学習



理論的にはこの二つはある種のトレードオフの関係にある。

## 推測

(より数理統計的)



↔ 肺癌

第〇〇遺伝子が肺癌に寄与  
有意水準5%

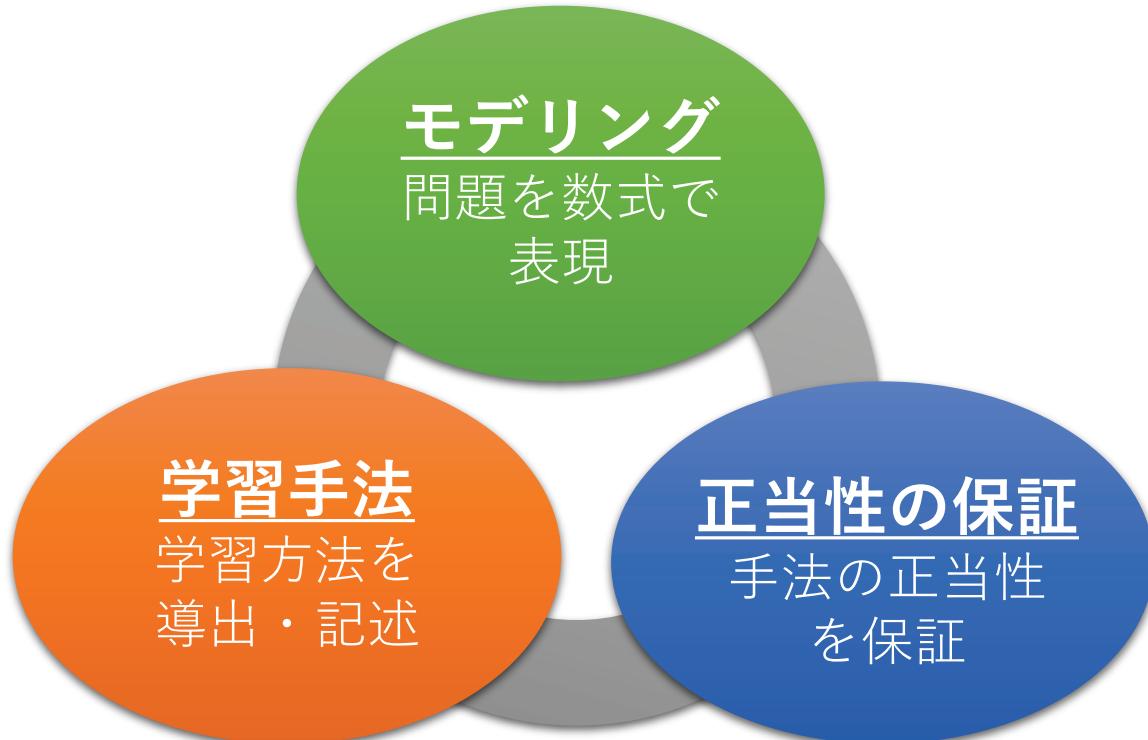
- 原因の究明.
- 仮説検定は典型例.

例：線形回帰分析



# 機械学習における数学の役割

1. 問題の数学的定式化（モデリング）
2. 手法の導出と記述（アルゴリズム）
3. 正当性の保証（学習理論）



記述言語：数学

- 確率-統計, 線形代数, 関数解析, 最適化理論

# 機械学習の歴史

# 機械学習と人工知能の歴史

1946: ENIAC, 高い計算能力

フォン・ノイマン「俺の次に頭の良い奴ができた」



1952: A. Samuelによるチェックカーズプログラム

統計的学習

ルールベース

1960年代前半:  
ELIZA(イライザ),  
擬似心理療法士

1980年代:  
エキスパートシステム

人手による学習ルール  
の作りこみの限界  
「膨大な数の例外」

Siriなどにつながる

1957: Perceptron, ニューラルネットワークの先駆け

第一次ニューラルネットワークブーム

1963: 線形サポートベクトルマシン

線形モデルの限界

1980年代: 多層パーセプトロン, 誤差逆伝搬,  
畳み込みネット

第二次ニューラルネットワークブーム

1992: 非線形サポートベクトルマシン  
(カーネル法)

1996: スパース学習 (Lasso)

2003: トピックモデル (LDA)

2012: Supervision (Alex-net)

非凸性の問題

データの増加  
+ 計算機の強化

第三次ニューラルネットワークブーム

人間

# 四則演算 単純なルール

機械

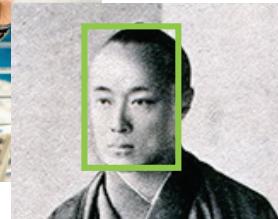
難  $193707721 \times 761838257287 - 2^{67} = -1$  易

人の手でプログラムするのは無理  
learn without being explicitly programmed

技術より

ノウハウ

易



人によって顔が違う、照明の当たり方で見え方・色が変わる、  
表情の違い、髪型の違い、顔の向きの違い、. . . .

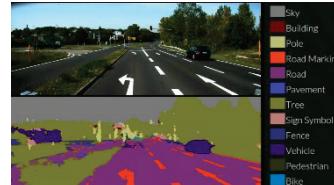
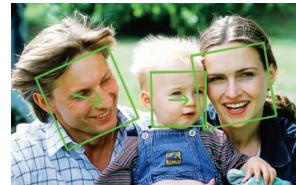
# 統計的学習の考え方

- 人がプログラムするのは認識の仕方ではなく学習の仕方

→数学で記述

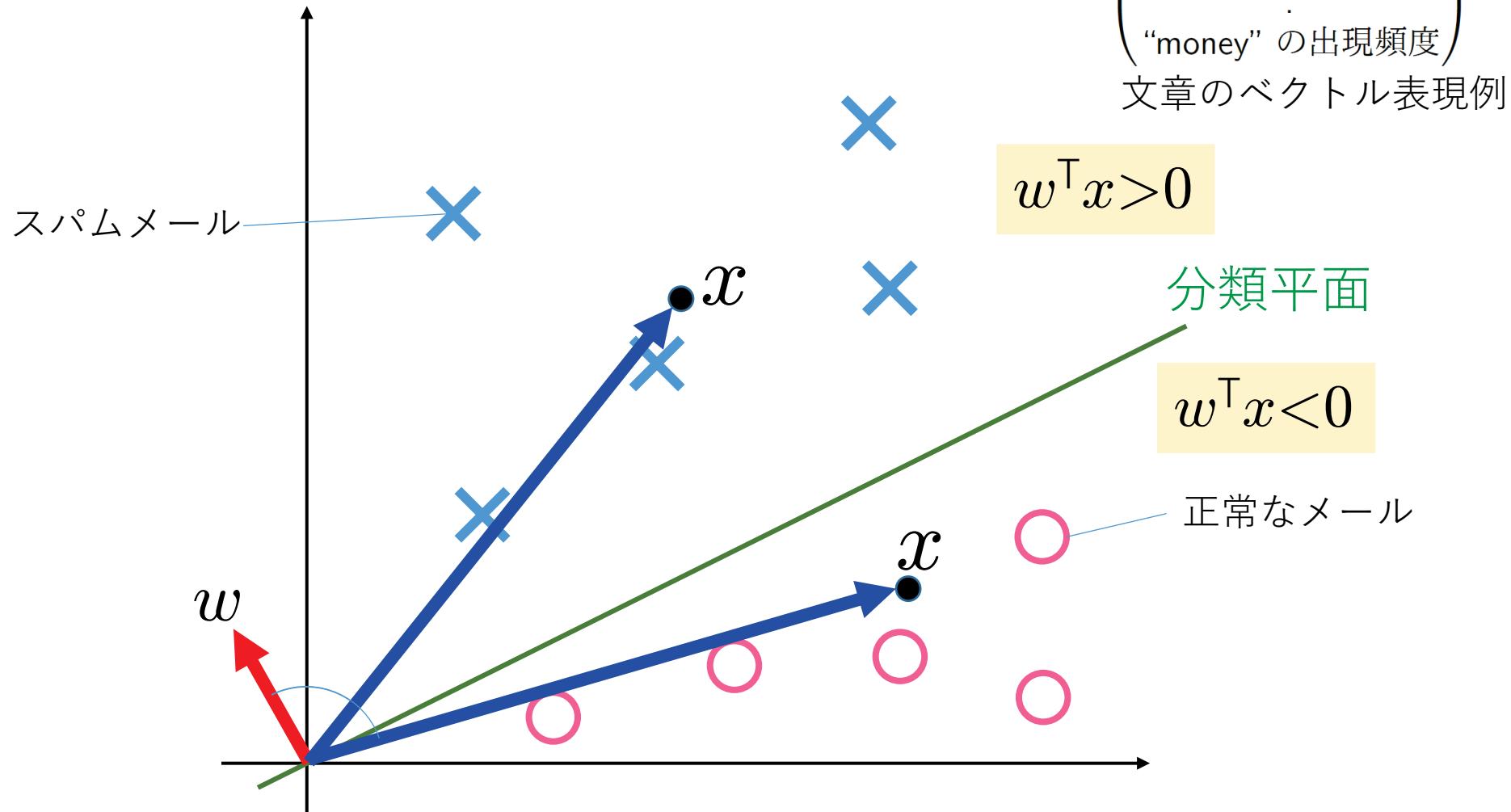


- 強い将棋ソフトを作りたい → 大量の棋譜データで学習
- 顔認識ソフトを作りたい → 大量の画像データで学習
- 車道を認識したい → 大量の車載カメラ画像で学習



# 線形分類機

$$w^\top x = w_1x_1 + w_2x_2 + \cdots + w_dx_d$$



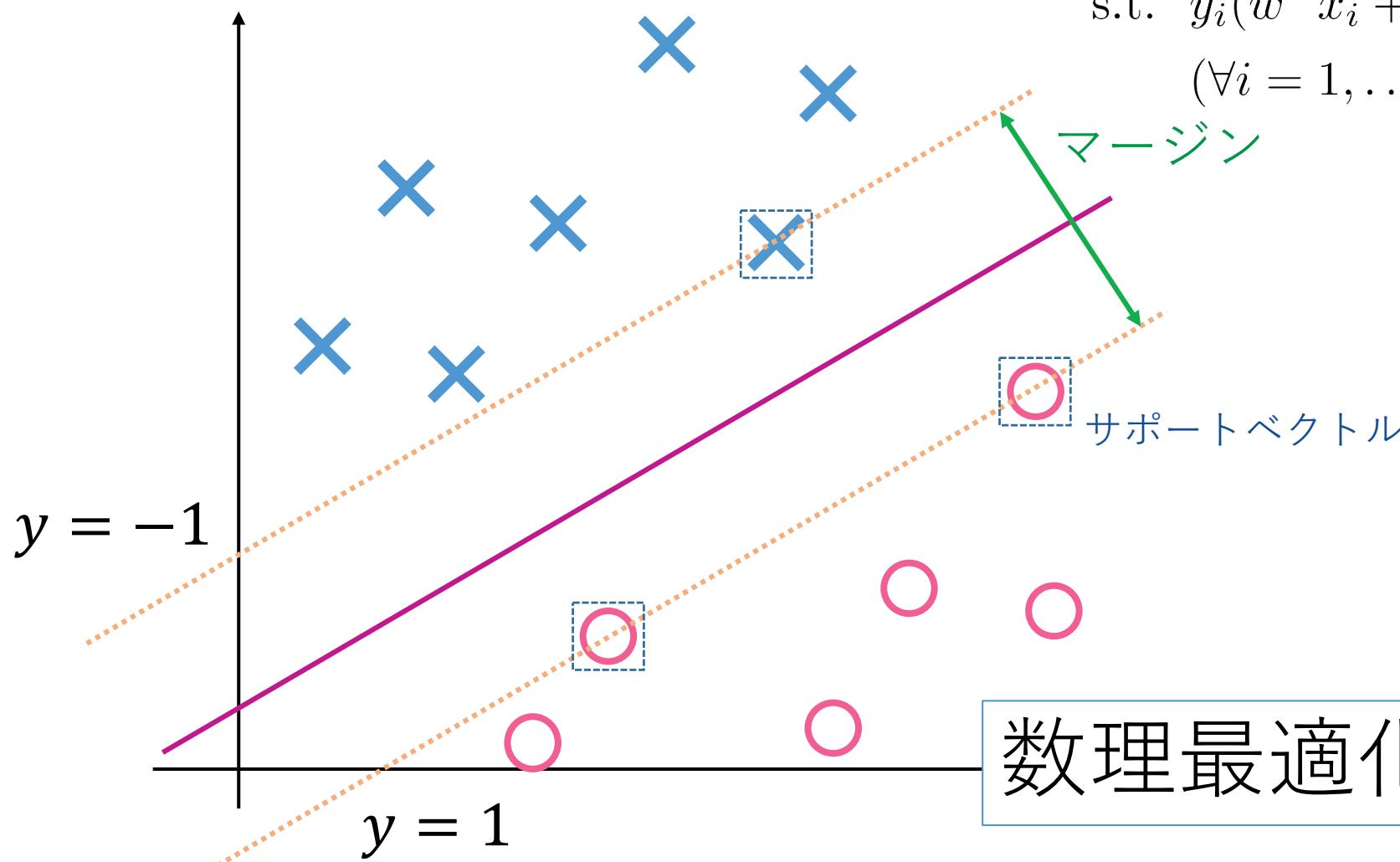
# サポートベクトルマシン (SVM)

[Vapnik,63]

マージンを最大化

VC (Vapnik-Chervonenkis) 理論による正当化

$$\begin{aligned} & \min_{w,b} \frac{\|w\|^2}{2} \\ \text{s.t. } & y_i(w^\top x_i + b) \geq 1 \\ & (\forall i = 1, \dots, n) \end{aligned}$$

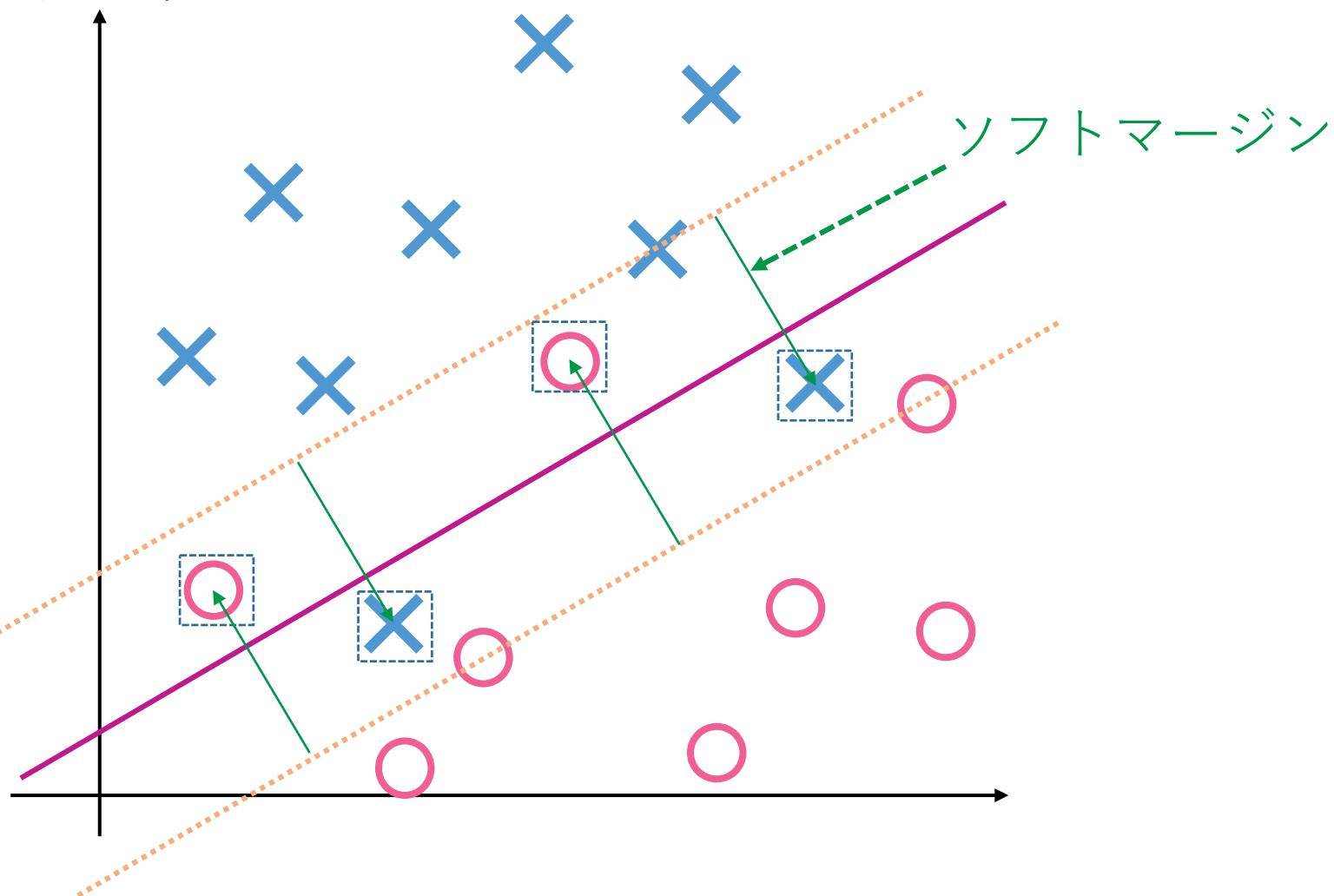


# ソフトマージンSVM

[Cortes+Vapnik,95]

マージンを最大化  
誤分類も許す

$$\min_{w,b} \sum_{i=1}^n \max\{1 - y_i(w^\top x_i + b), 0\} + C \frac{\|w\|^2}{2}$$



# 機械学習と人工知能の歴史

1946: ENIAC, 高い計算能力

フォン・ノイマン「俺の次に頭の良い奴ができた」

1952: A. Samuelによるチェックカーズプログラム

統計的学習

1957: Perceptron, ニューラルネットワークの先駆け

第一次ニューラルネットワークブーム

1963: 線形サポートベクトルマシン

線形モデルの限界

1980年代: 多層パーセプトロン, 誤差逆伝搬,

畳み込みネット

第二次ニューラルネットワークブーム

1992: 非線形サポートベクトルマシン

(カーネル法)

1996: スパース学習 (Lasso)

非凸性の問題

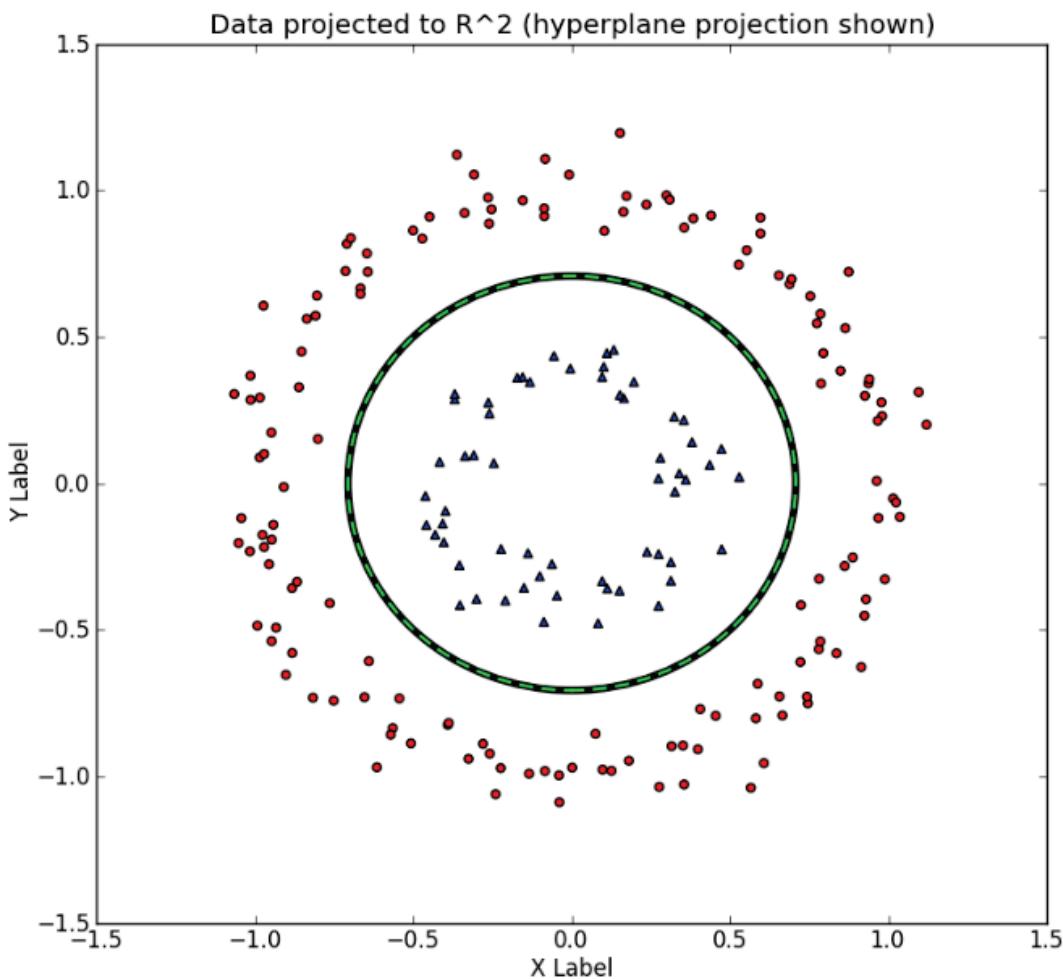
2003: トピックモデル (LDA)

データの増加  
+ 計算機の強化

2012: Supervision (Alex-net)

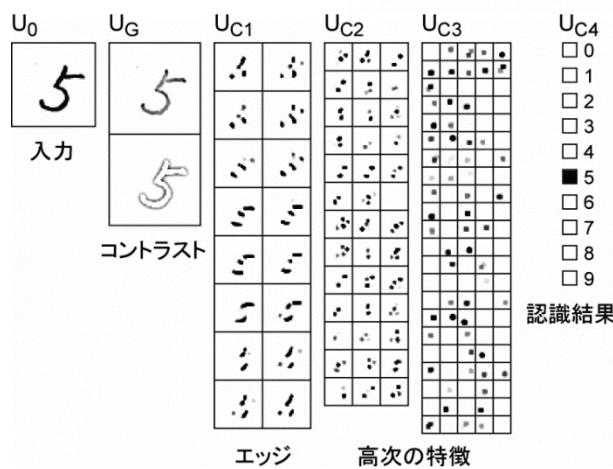
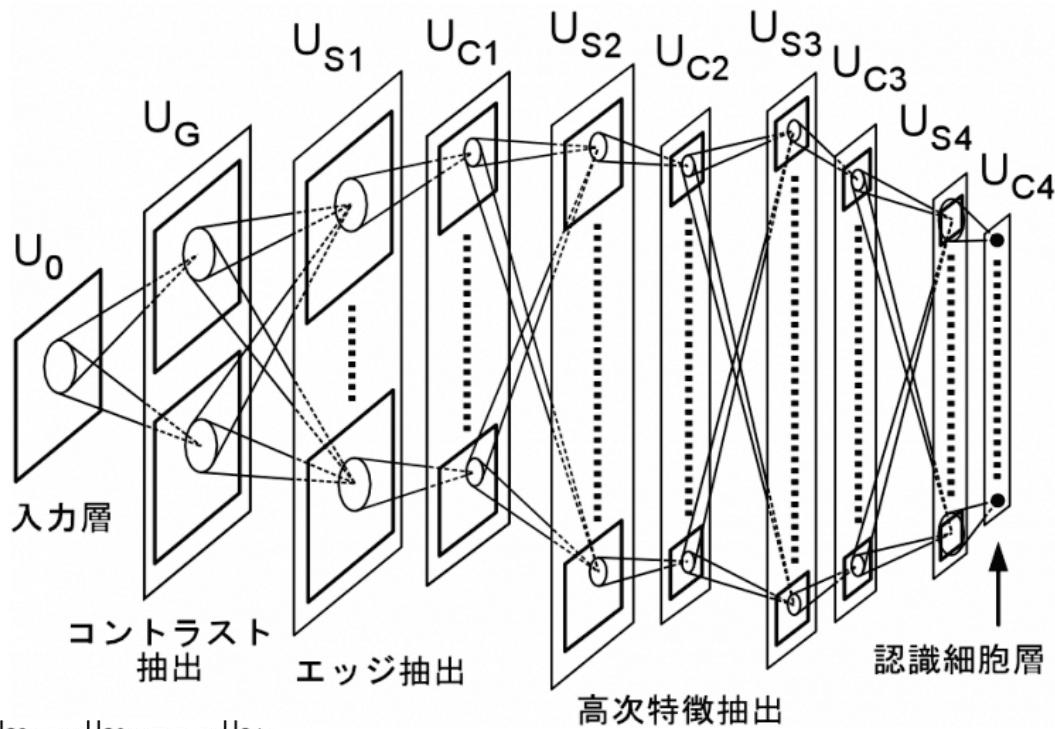
第三次ニューラルネットワークブーム

# 非線形判別



# ネオコグニトロン

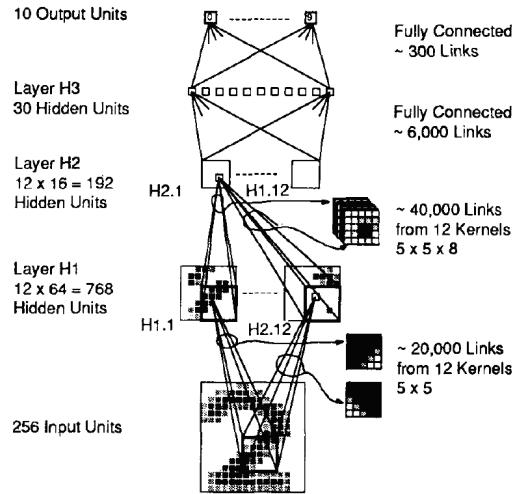
[福島, 79]



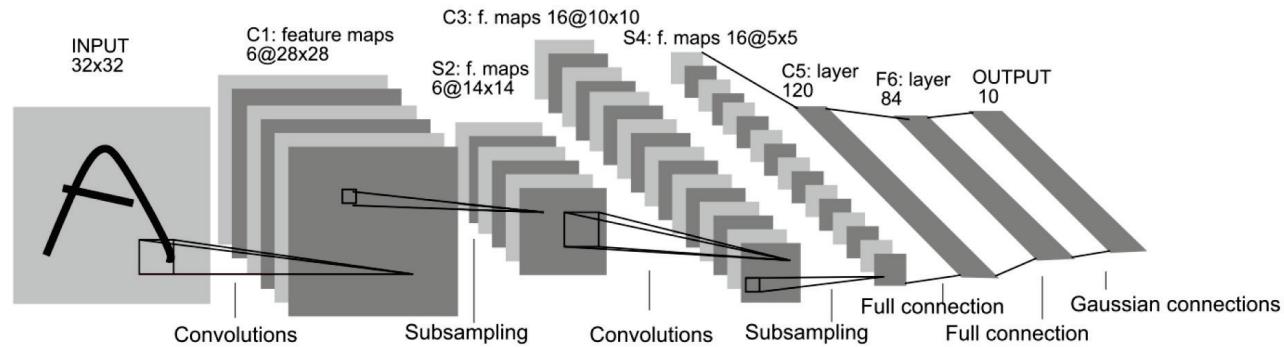
- ・人間の脳を模倣
- ・畳み込みネットの初期型
- ・自己組織型学習  
→素子を足していく

# LeNet

[LeCun+etal,89]



LeNet-5  
[LeCun et al,98]



- 畳み込み + プーリング：現在も使われている構造
- 誤差逆伝搬法でパラメータを更新
- 手書き文字認識データセット（MNIST）で99%の精度を達成

# 機械学習と人工知能の歴史

1946: ENIAC, 高い計算能力

フォン・ノイマン「俺の次に頭の良い奴ができた」



1952: A. Samuelによるチェックカーズプログラム

統計的学習

ルールベース

1960年代前半:  
ELIZA(イライザ),  
擬似心理療法士

1980年代:  
エキスパートシステム

人手による学習ルール  
の作りこみの限界  
「膨大な数の例外」

Siriなどにつながる

1957: Perceptron, ニューラルネットワークの先駆け

第一次ニューラルネットワークブーム

1963: 線形サポートベクトルマシン

線形モデルの限界

1980年代: 多層パーセプトロン, 誤差逆伝搬,  
畳み込みネット

第二次ニューラルネットワークブーム

1992: 非線形サポートベクトルマシン  
(カーネル法)

1996: スパース学習 (Lasso)

2003: トピックモデル (LDA)

2012: Supervision (Alex-net)

非凸性の問題

データの増加  
+ 計算機の強化

第三次ニューラルネットワークブーム

# 問題点

- 誰でも実装できるわけではなかった。  
e.g. 「LeNetはYan LeCunしか実装できない」といった噂
- 様々な職人芸的なノウハウが存在。
  - パラメータのチューニング：学習率，層の数，層の幅
- 大域的最適解の計算が難しい。  
局所最適解しか得られない。

→これらは現代でも未解決  
(実装に関してはライブラリの充実でかなり解決)

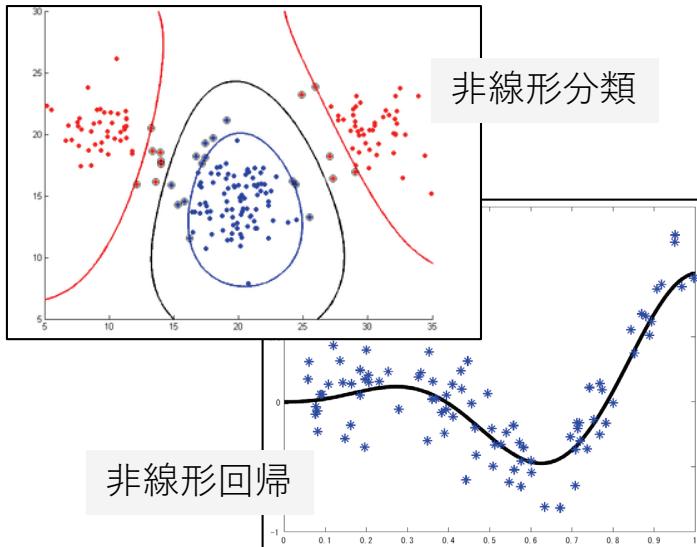
誰でも実装できて、最適解が一つの手法が欲しい。  
→ カーネルを用いたサポートベクトルマシン

# カーネル法

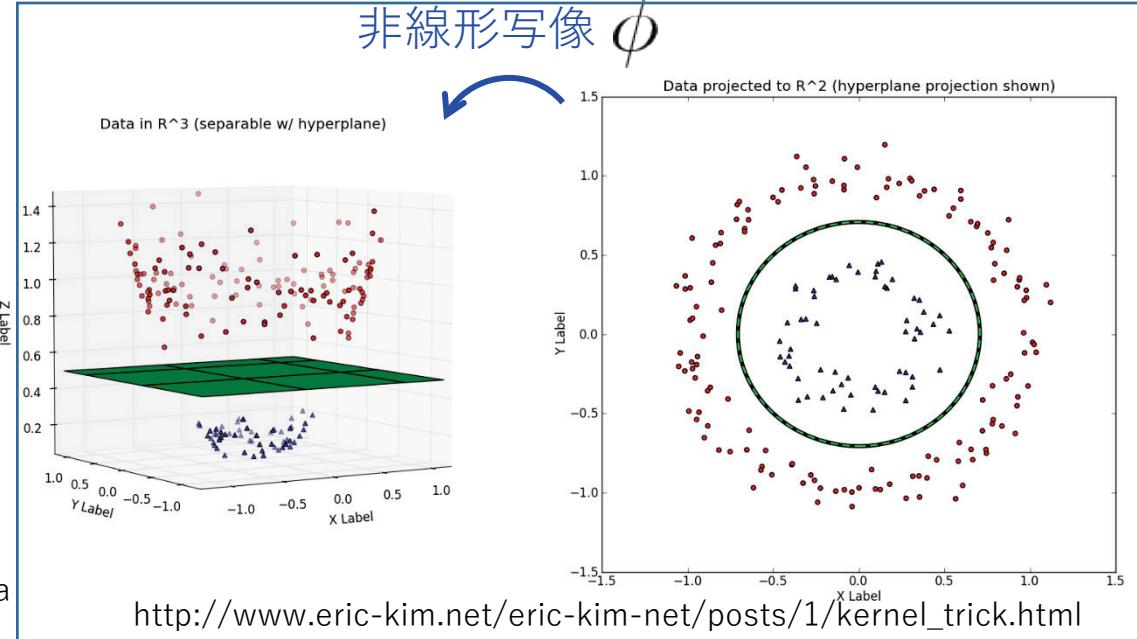
$$\min_{\alpha_i, b} \sum_{i=1}^n \max \left\{ 1 - y_i \left( \sum_{j=1}^n k(x_j, x_i) \alpha_j + b \right), 0 \right\} + C \sum_{i,j=1}^n \alpha_i \alpha_j k(x_i, x_j)$$

カーネルトリック

$$k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$$

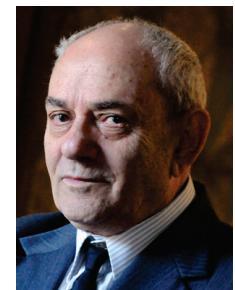


<http://wiki.eigenvalue.com/index.php?title=Svmda>



## 関数解析：再生核ヒルベルト空間の理論

- 凸最適化問題で解ける。
  - ✓ 効率的な最適化手法が存在。
  - ✓ 解は一つ。誰が解いても同じ答えが返ってくる。
- VC理論・経験過程の理論による汎化誤差の保証。



Vladimir Vapnik

$$\|\hat{f} - f_0\|_{L_2}^2 \leq O_p(n^{-\frac{1}{1+s}})$$

# 90年代以降

## データ解析としての機械学習

統計学やデータマイニングとの融合

- 高次元スペース学習
- ベイズモデリング
- オンライン学習, 確率的最適化  
→ ビッグデータ解析への活用

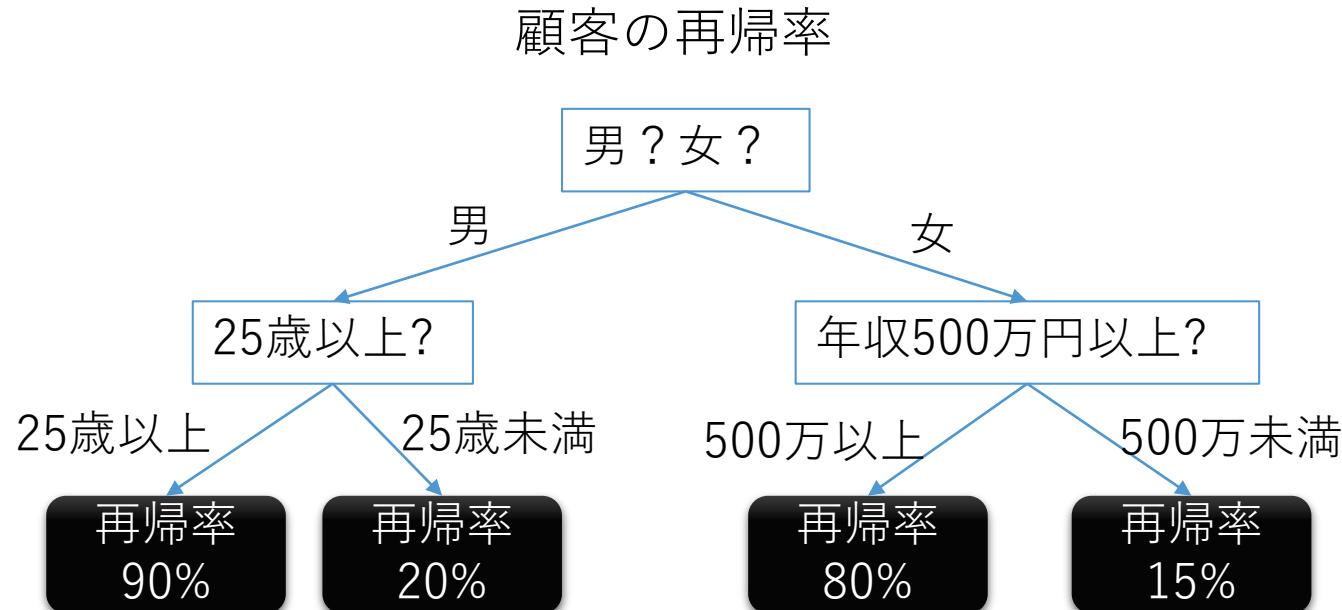


メディアに出る「人工知能」はここに属することが多い

- データの増加 + 計算機の強化  
→ 深層学習再興, **第三次ニューラルネットワークブーム**へ

# 決定木

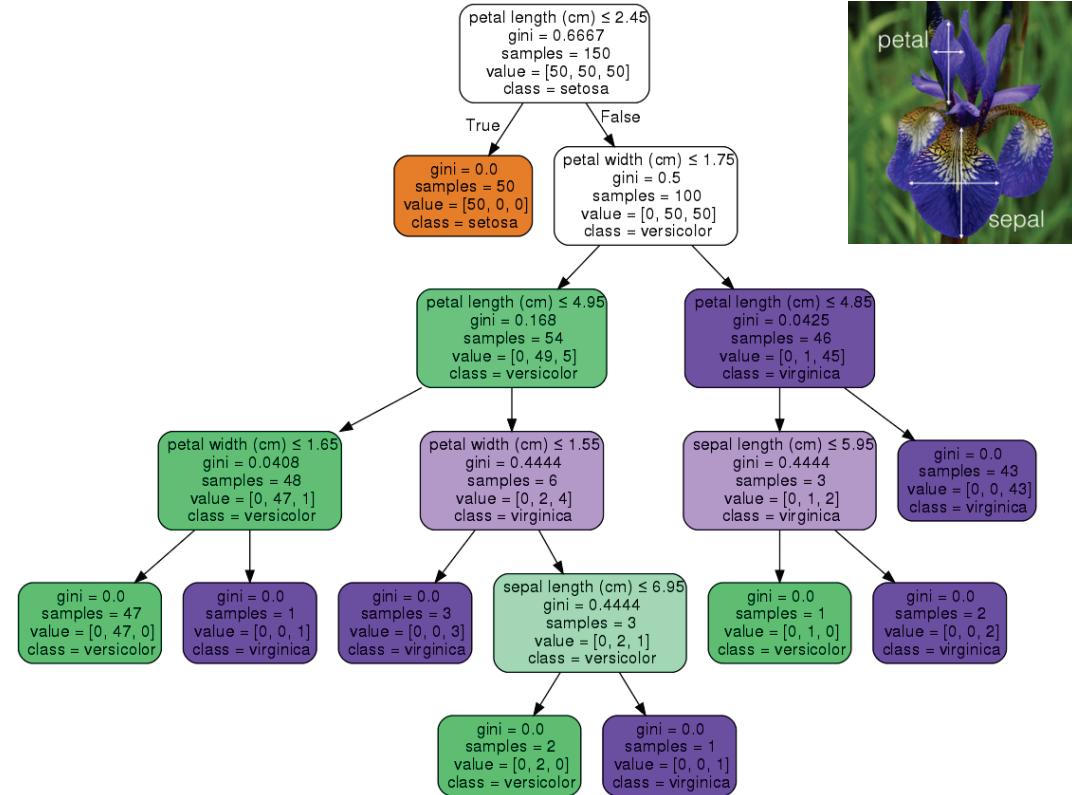
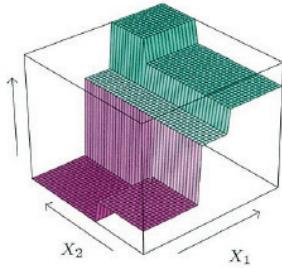
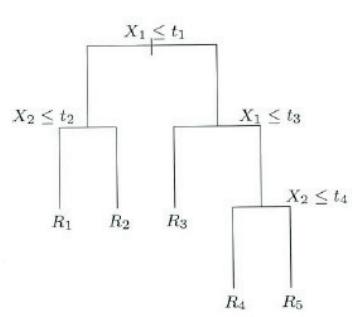
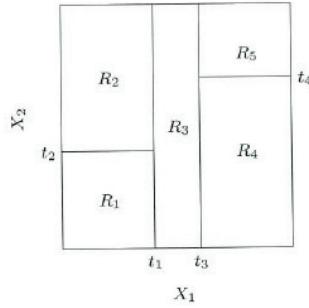
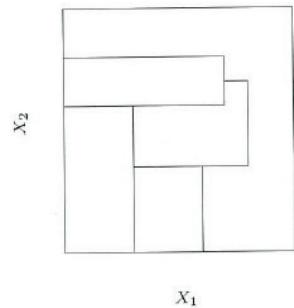
- ・解釈可能性が高い
- ・決定木を組み合わせた勾配ブースティングはデータマイニング系コンペティションで常連  
(判別だけでなく回帰にも利用可能)



- ・学習された決定木から要因を把握しやすい。
- ・分析結果から対策を立てるのに有用。

# 決定木

## 決定木の様子 2次元の説明変数

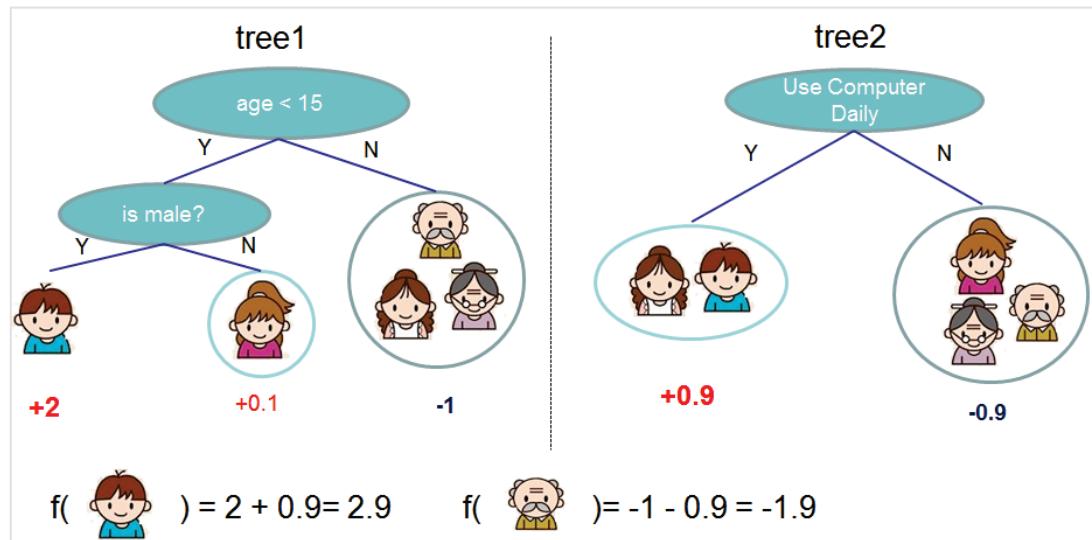


図はHastie, Tibshirani, Friedman: The Elements of Statistical Learning, Springer, 2001.より

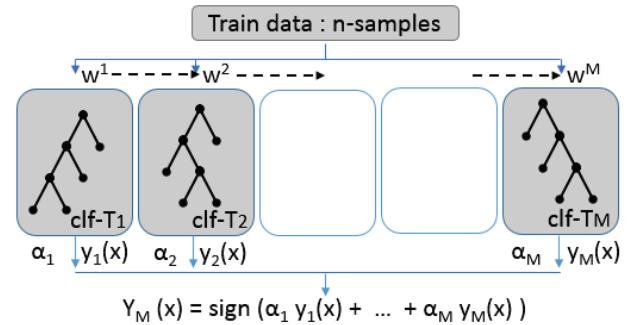
Python scikit-learnによるiris(アヤメ)データ分類

# 勾配ブースティング

- XGBoostやLightGBMが有名
- 「決定木の和」で強力な判別を実現
  - 決定木一つでは複雑な判別が難しい→ 複数用意してその和(多数決)を取る
  - 和の取り方に勾配ブースティングと呼ばれる技法を使用



「コンピュータゲームが好きか？」を判別



沢山の決定木の多数決を出力

[Chen, Guestrin: XGBoost: A Scalable Tree Boosting System. KDD2016.]

[Ke, Meng, Finley, Wang, Chen, Ma, Ye, Liu: LightGBM: A Highly Efficient Gradient Boosting Decision Tree. NIPS2017.]

# 勾配ブースティング

XGBoostは各種データ解析コンペティションで好成績

## Machine Learning Challenge Winning Solutions

XGBoost is extensively used by machine learning practitioners to create state of art data science solutions, this is a list of machine learning winning solutions with XGBoost. Please send pull requests if you find ones that are missing here.

- Maksims Volkovs, Guangwei Yu and Tomi Poutanen, 1st place of the [2017 ACM RecSys challenge](#). Link to [paper](#).
- Vlad Sandulescu, Mihai Chiru, 1st place of the [KDD Cup 2016 competition](#). Link to [the arxiv paper](#).
- Marios Michailidis, Mathias Müller and HJ van Veen, 1st place of the [Dato Truly Native? competition](#). Link to [the Kaggle interview](#).
- Vlad Mironov, Alexander Guschin, 1st place of the [CERN LHCb experiment Flavour of Physics competition](#). Link to [the Kaggle interview](#).
- Josef Slavicek, 3rd place of the [CERN LHCb experiment Flavour of Physics competition](#). Link to [the Kaggle interview](#).
- Mario Filho, Josef Feigl, Lucas, Gilberto, 1st place of the [Caterpillar Tube Pricing competition](#). Link to [the Kaggle interview](#).
- Qingchen Wang, 1st place of the [Liberty Mutual Property Inspection](#). Link to [the Kaggle interview](#).
- Chenglong Chen, 1st place of the [Crowdflower Search Results Relevance](#). Link to [the winning solution](#).
- Alexandre Barachant ("Cat") and Rafał Cycoń ("Dog"), 1st place of the [Grasp-and-Lift EEG Detection](#). Link to [the Kaggle interview](#).
- Halla Yang, 2nd place of the [Recruit Coupon Purchase Prediction Challenge](#). Link to [the Kaggle interview](#).
- Owen Zhang, 1st place of the [Avito Context Ad Clicks competition](#). Link to [the Kaggle interview](#).
- Keiichi Kuroyanagi, 2nd place of the [Airbnb New User Bookings](#). Link to [the Kaggle interview](#).
- Marios Michailidis, Mathias Müller and Ning Situ, 1st place [Homesite Quote Conversion](#). Link to [the Kaggle interview](#).

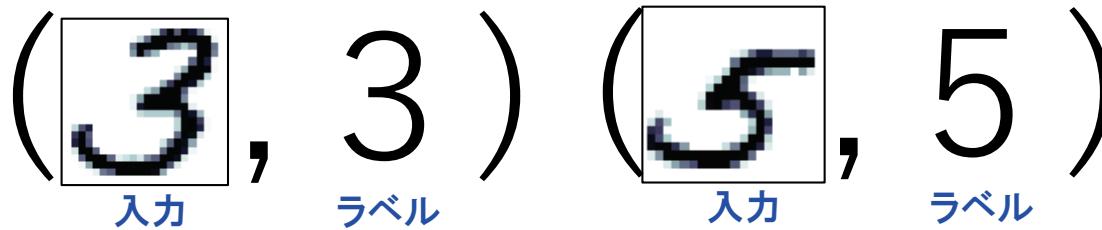
# 機械学習の数理

# 機械学習の問題設定

教師あり学習：

データ： $(x, y) \leftarrow$ ある入力 $x$ とそれに対するラベル $y$ の組

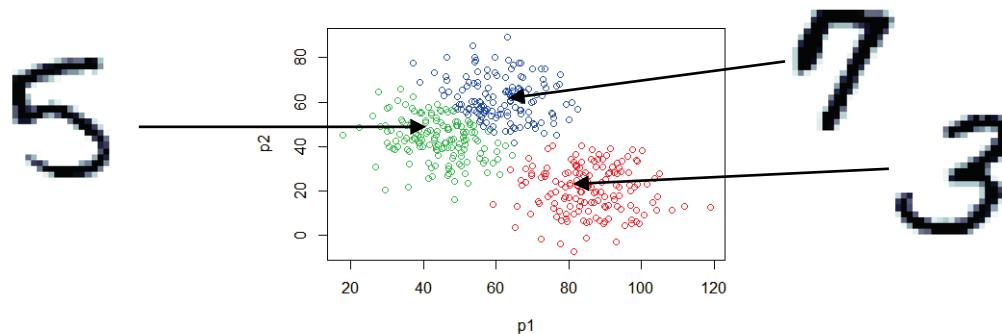
問題の例：回帰，判別



教師なし学習：

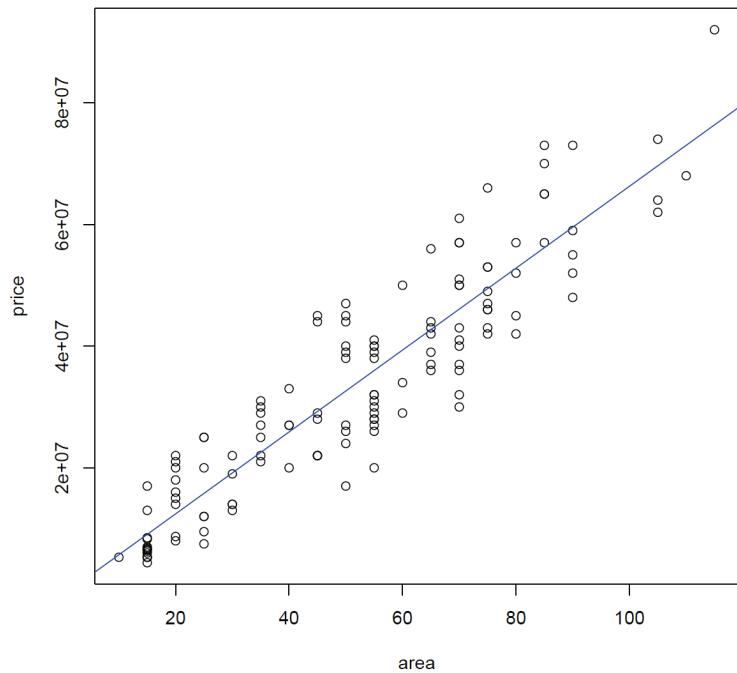
データ： $(x) \leftarrow$ ラベルがない

問題の例：クラスタリング，音源分離，異常検知

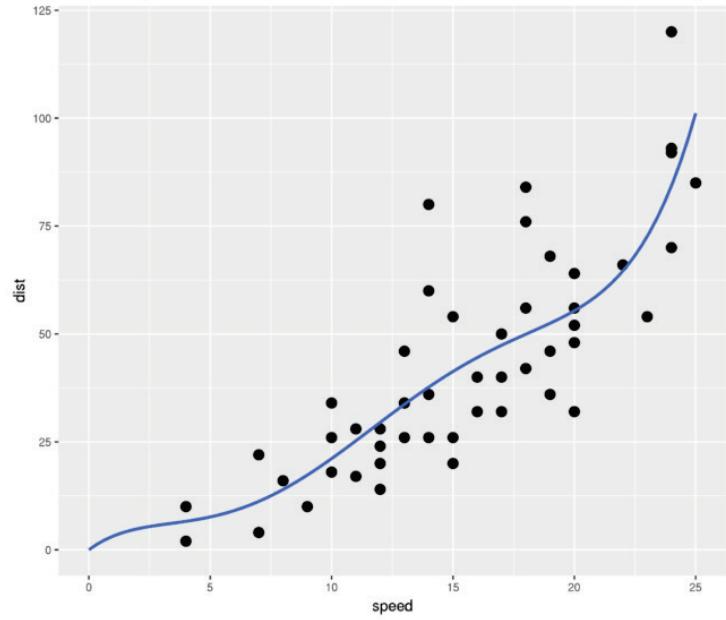


半教師有り学習：ラベルの付いているデータと付いてないデータが混在

# 回帰



線形

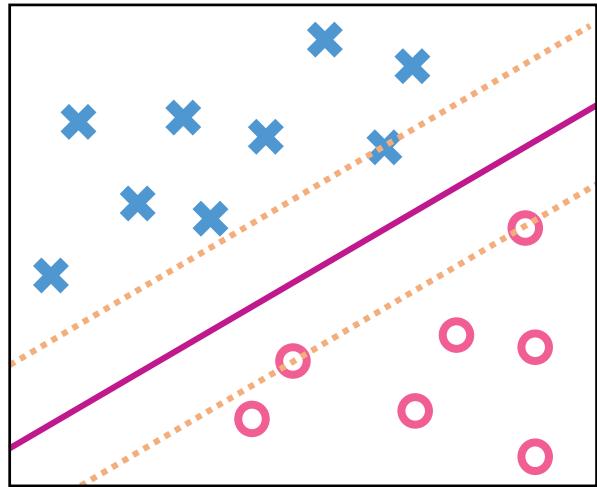


非線形

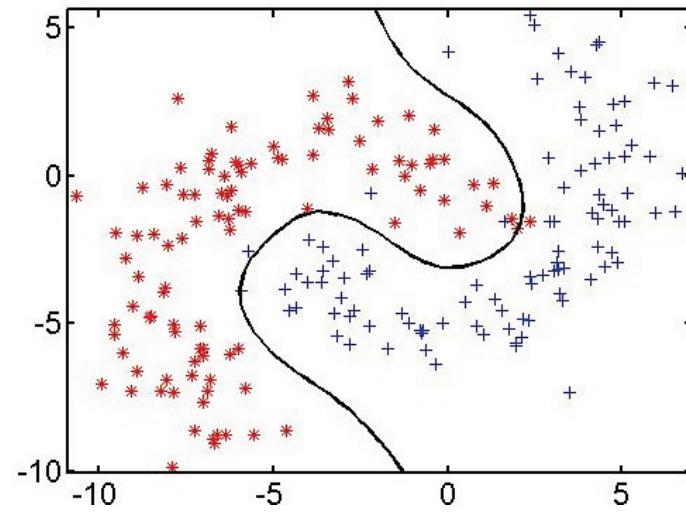
入力xから実数の出力yを予測

- 線形
- 非線形

# 判別



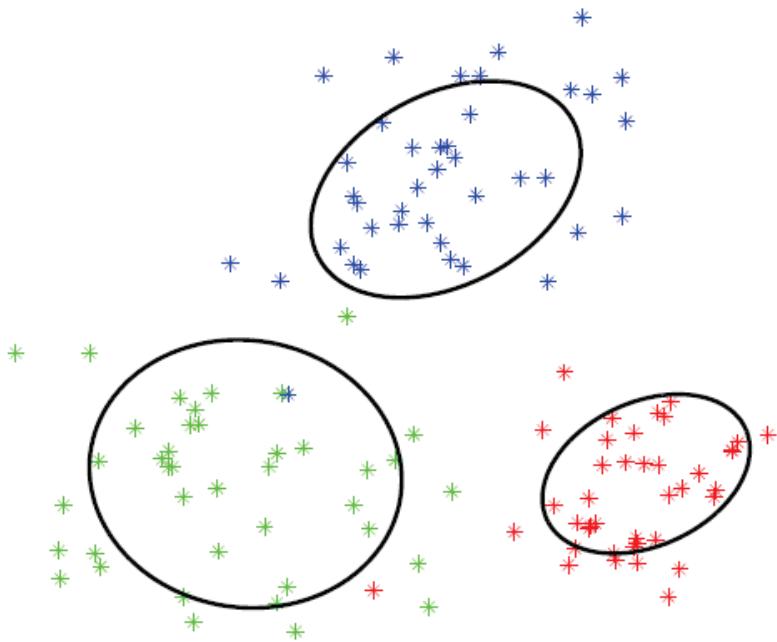
線形



非線形

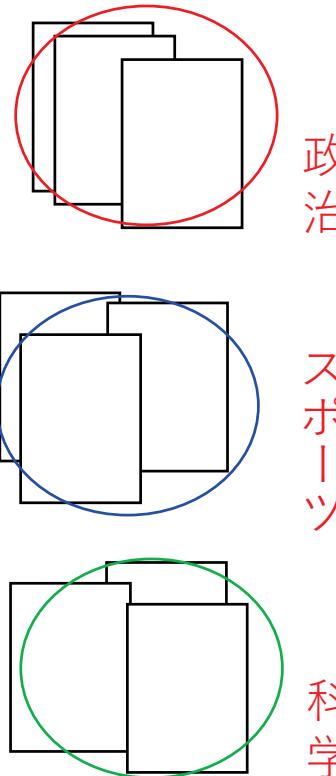
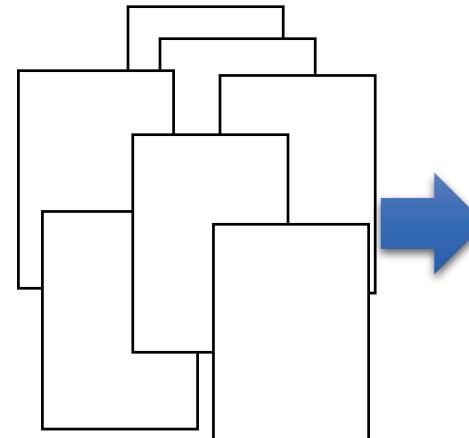
- 入力 $x$ からカテゴリーの出力 $y$ を予測
- 線形
  - 非線形

# クラスタリング



混合ガウス分布によるクラスタリング

文章データ  
(ニュース記事など)

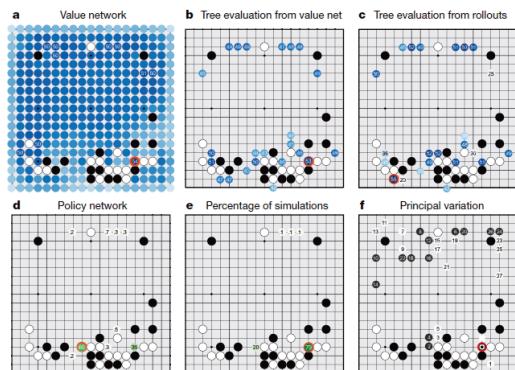


トピックモデル  
• 文章分類

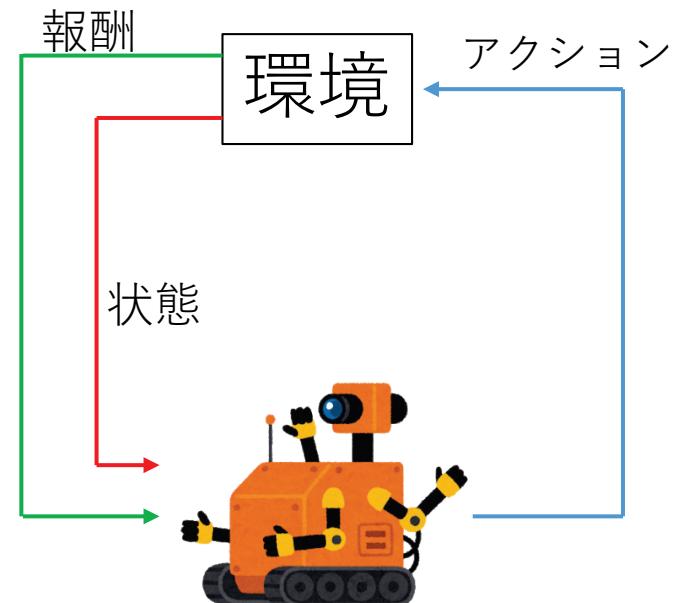
# 強化学習



Google research blog, 8/March/2016.  
“Deep Learning for Robots: Learning from Large-Scale Interaction.”



[Silver et al. (Google Deep Mind): Mastering the game of Go with deep neural networks and tree search,  
Nature, 529, 484—489, 2016]

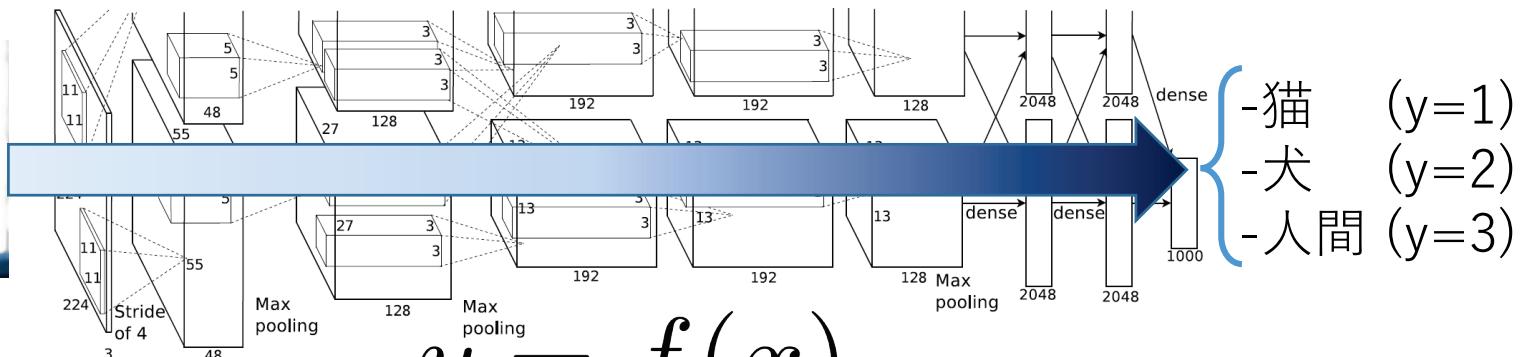


# 教師あり学習

画像

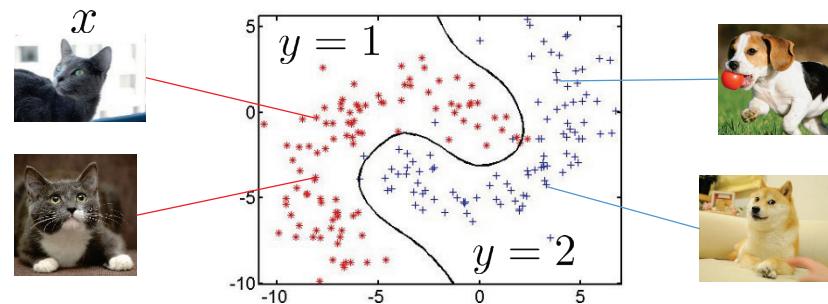


$x$



$$y = f(x)$$

$y$



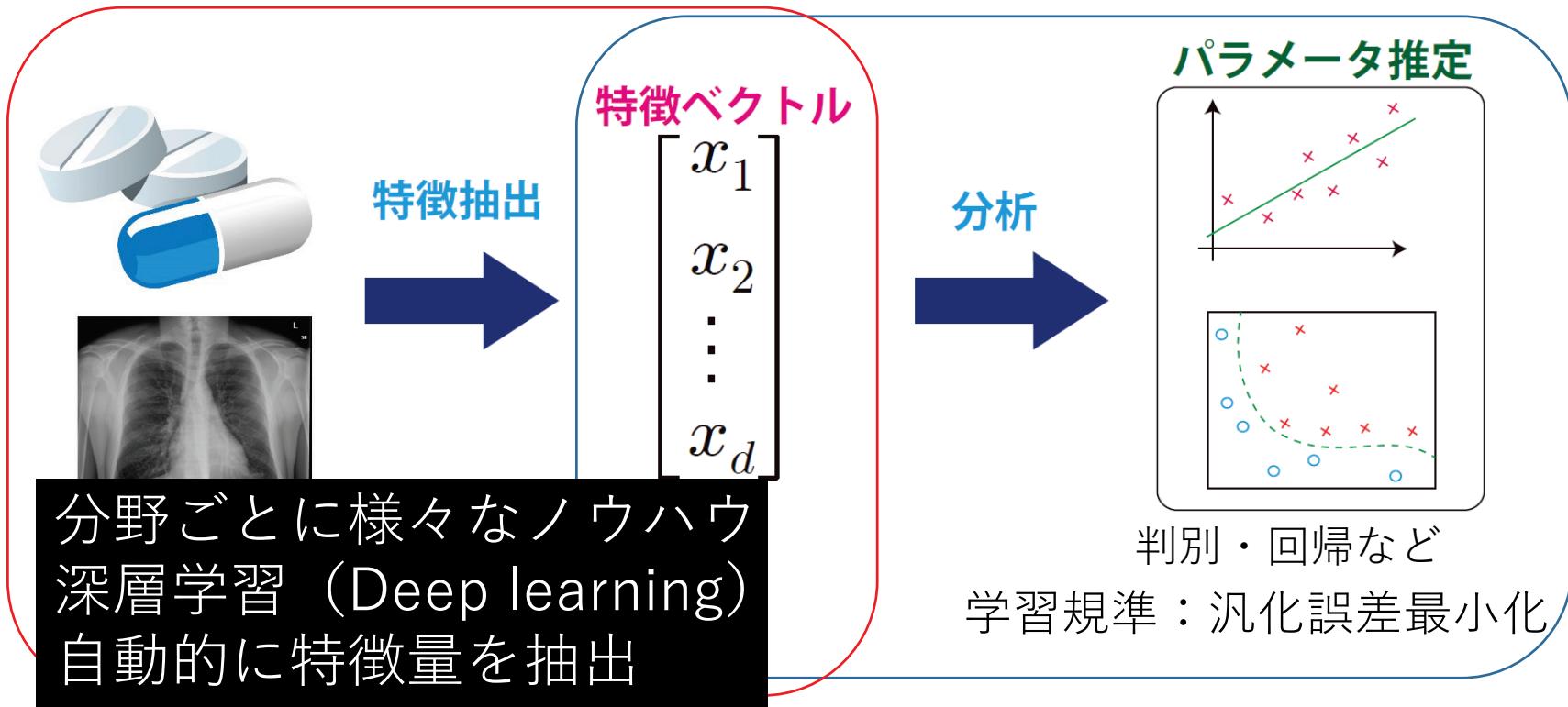
学習：「関数」をデータに当てはめる

モデル：関数の集合（例：深層NNの表せる関数の集合）

# 予測モデルの学習

問題ごと

共有化可能



予測モデルの構築

$$y = f(x; \theta)$$

モデルの  
パラメータ

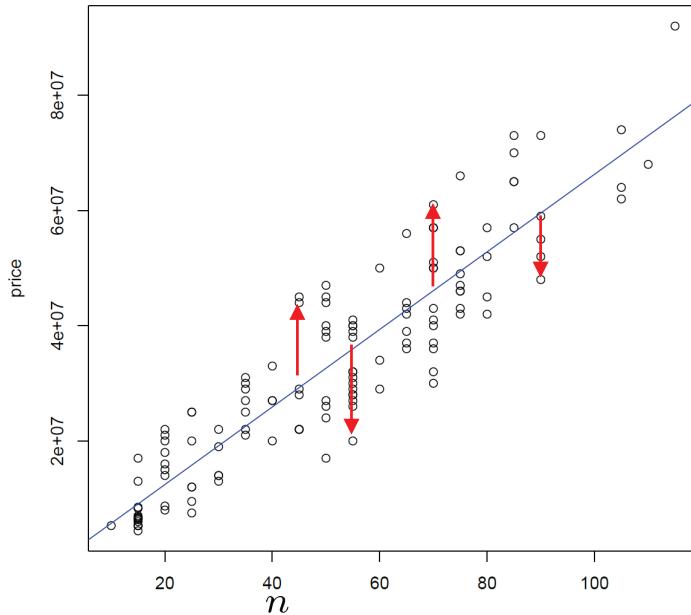
一度特徴ベクトルに変換してしまえばあとは統計の問題。  
→汎用的な手法（機械学習）を適用できる。

# 線形モデル

$$y = \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p + \beta_0 + \epsilon$$

$y$ :従属変数,  $x$ :特徴ベクトル

マンション価格 =  $\beta_1 \times$  床面積 +  $\beta_2 \times$  築年数 +  $\beta_3$  + (揺らぎ)



最小二乗法

$$\min_{\beta_0, \beta_1, \beta_2, \beta_3} \sum_{i=1}^n (y_i - \beta_1 x_{i,1} - \beta_2 x_{i,2} - \beta_3 x_{i,3} - \beta_0)^2$$

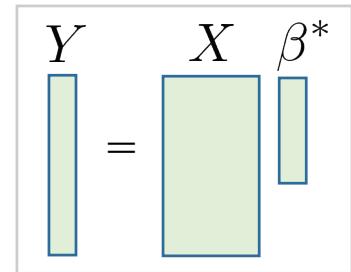
# 最小二乗法

$n$ 個の観測値（サンプル）： $(y_i, \mathbf{x}_i) \in \mathbb{R} \times \mathbb{R}^d$  ( $i = 1, \dots, n$ )

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n, \quad X = \begin{bmatrix} \mathbf{x}_1^\top & 1 \\ \vdots & \vdots \\ \mathbf{x}_n^\top & 1 \end{bmatrix} \in \mathbb{R}^{n \times (d+1)}, \quad \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix} \in \mathbb{R}^n$$

$\beta^*$ を真の回帰係数（これを推定したい）とすると、

$$Y = X\beta^* + \epsilon$$

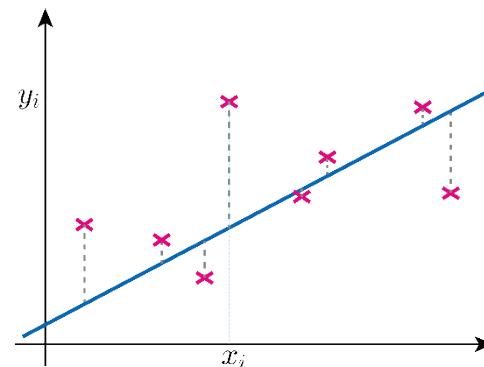


**最小二乗推定量（最尤推定量）：**

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^{d+1}} \sum_{i=1}^n (y_i - [\mathbf{x}_i^\top 1] \beta)^2$$

$$= \arg \min_{\beta \in \mathbb{R}^{d+1}} \|Y - X\beta\|^2$$

$$= (X^\top X)^{-1} X^\top Y$$



# 訓練誤差と汎化誤差

パラメータ  $\theta$  : データの構造を表す変数 (例: 判別平面)

損失関数  $\ell(Y, f(X, \theta))$  : パラメータ  $\theta$  がデータをどれだけ説明しているか

**汎化誤差** : 損失の期待値

$$\mathbb{E}[\ell(Y, f(X, \theta))]$$

本当は最小化したいもの.

※クラスタリング等, 教師なし学習も尤度を使ってこのように書ける.

**訓練誤差** : 有限個のデータで代用

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i, \theta))$$

代わりに最小化するもの.

この二つには大きなギャップがある.  
[過学習]

# 基本的な考え方

- パラメータを  $\theta$  としたときに、観測されたデータが観測される確率（尤度）

尤度 
$$\prod_{i=1}^n p(z_i|\theta) \quad : \text{確率モデル}$$

尤度が高ければ、観測データが観測される確率が高い → 「尤もらしい」

負の対数尤度 
$$\sum_{i=1}^n -\log(p(z_i|\theta))$$
  

$$\ell(z_i, \theta)$$

→ 最小化で観測データを良く表現するパラメータが得られる。

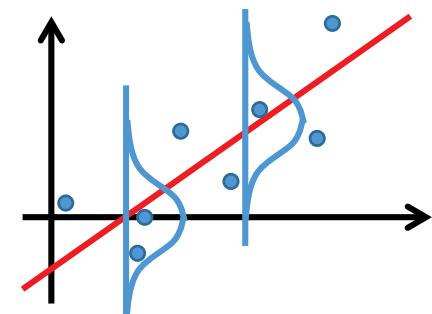
**「最尤推定」**

(ベイズ推定も重要だがここでは割愛)

## 線形回帰

$$p(y_i|x_i, \theta) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y_i - x_i^\top \theta)^2}{2}\right) \quad \begin{array}{l} \text{正規分布} \\ \text{平均 } x_i^\top \theta, \text{ 分散 } 1 \end{array}$$

$$-\log(p(y_i|x_i, \theta)) = \frac{(y_i - x_i^\top \theta)^2}{2} + C \rightarrow \text{最小二乗法}$$



# KL-divergence

$$D(p||q) = \int p(z) \log \left( \frac{p(z)}{q(z)} \right) dz$$

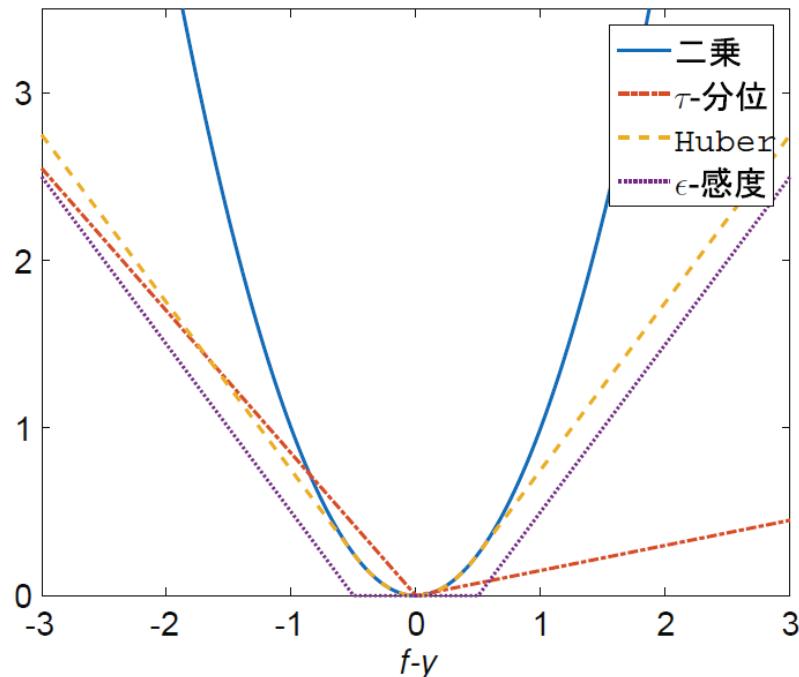
真の分布    モデルの分布

$$D(p^*||p_\theta) = \underbrace{\int -p^*(z) \log (p_\theta(z)) dz + \int p^*(z) \log (p^*(z)) dz}_{\text{サンプル平均で代用}} \\ - \frac{1}{n} \sum_{i=1}^n \log (p_\theta(z_i))$$

対数尤度最大化はKL-divergence最小化の近似ともみなせる

# 回帰の損失関数

- 二乗損失:  $\ell(y, f) = \frac{1}{2}(y - f)^2.$
- $\tau$ -分位点損失:  $\ell(y, f) = (1 - \tau) \max\{f - y, 0\} + \tau \max\{y - f, 0\}.$   
ただし,  $\tau \in (0, 1)$ . 分位点回帰に用いられる.
- $\epsilon$ -感度損失:  $\ell(y, f) = \max\{|y - f| - \epsilon, 0\},$   
ただし,  $\epsilon > 0$ . サポートベクトル回帰に用いられる.



※各損失関数は必ずしも確率モデルと対応するわけではない

# 判別

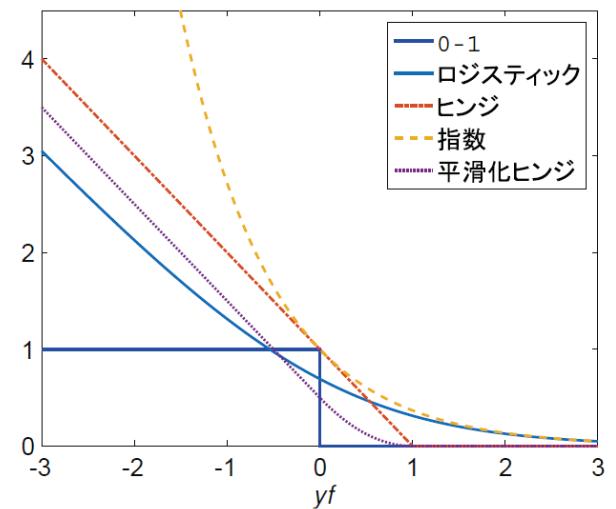
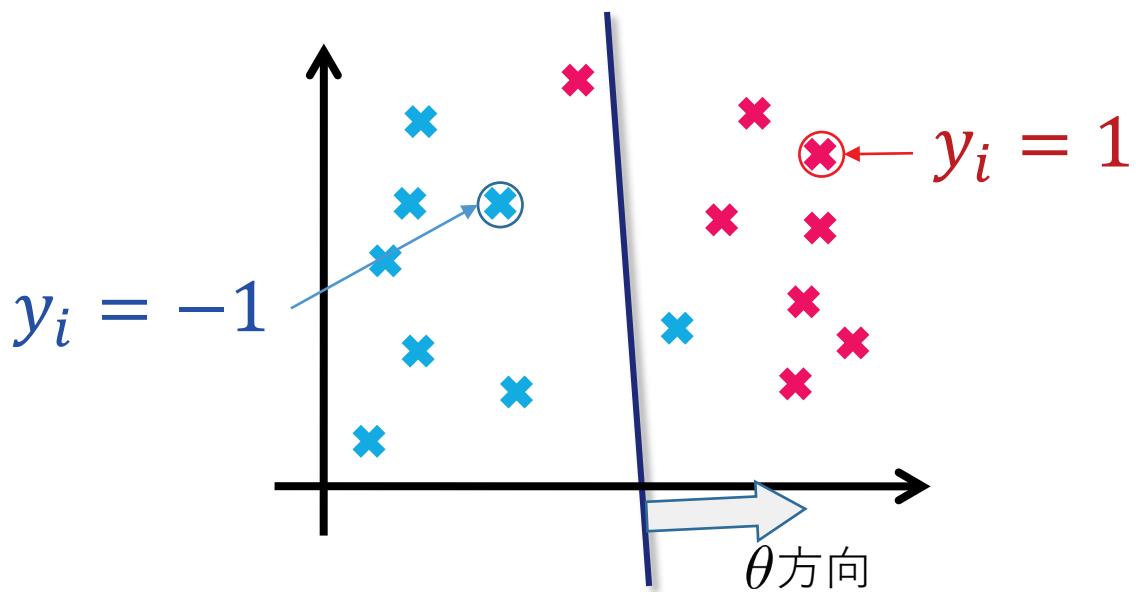
- 判別

## 損失関数

$$\ell(y, x^\top \theta) = \log(1 + \exp(-yx^\top \theta)) \quad (\text{ロジスティック損失})$$

## 訓練誤差最小化

$$\min_{\theta \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, x_i^\top \theta) = \min_{\theta \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i x_i^\top \theta)) \quad (\text{ロジスティック回帰})$$

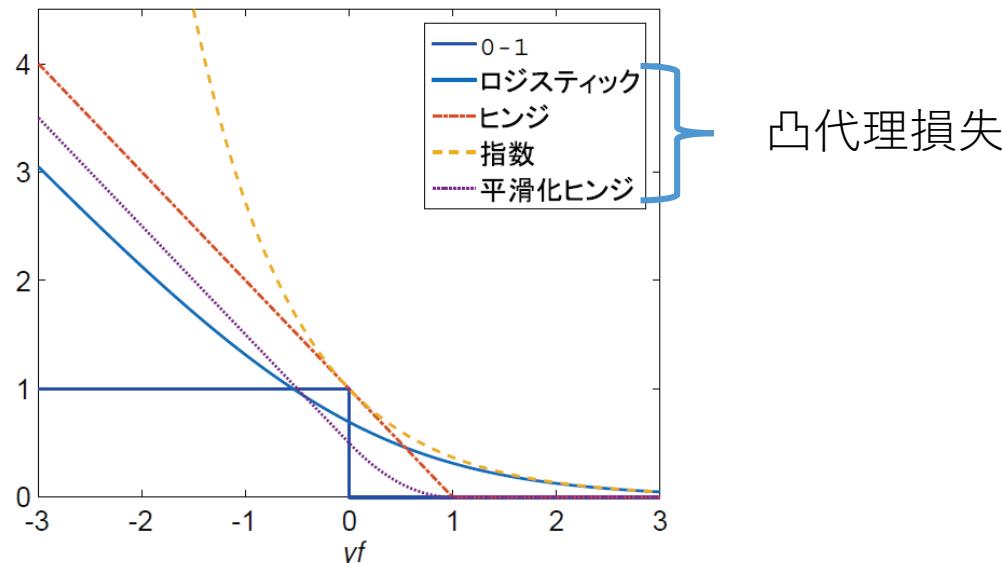


# 判別の損失関数

$$y \in \{\pm 1\}$$

- ロジスティック損失:  $\ell(y, f) = \log((1 + \exp(-yf))/2).$
- ヒンジ損失:  $\ell(y, f) = \max\{1 - yf, 0\}.$
- 指数損失:  $\ell(y, f) = \exp(-yf).$
- 平滑化ヒンジ損失:

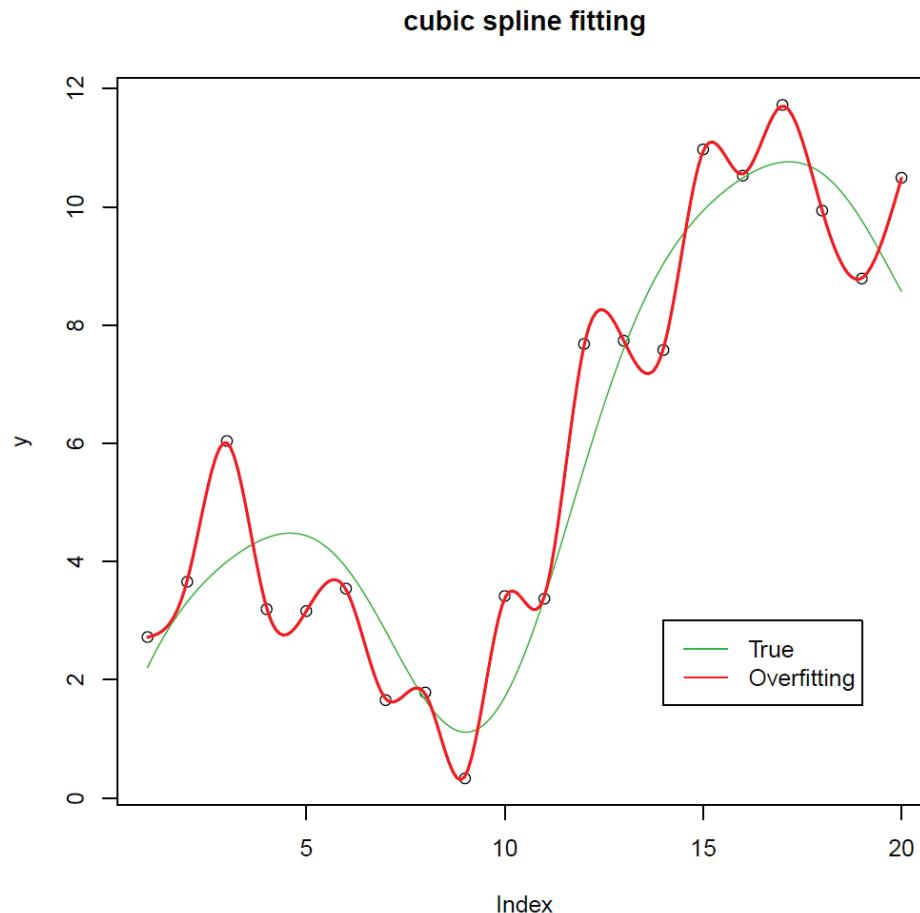
$$\ell(y, f) = \begin{cases} 0, & (yf \geq 1), \\ \frac{1}{2} - yf, & (yf < 0), \\ \frac{1}{2}(1 - yf)^2, & (\text{otherwise}). \end{cases}$$



# 過学習

複雑なモデル（例えば深層ニューラルネット）を用いるのが常に良い選択か？

→ そうとは限らない。 「過学習」に注意する必要あり。



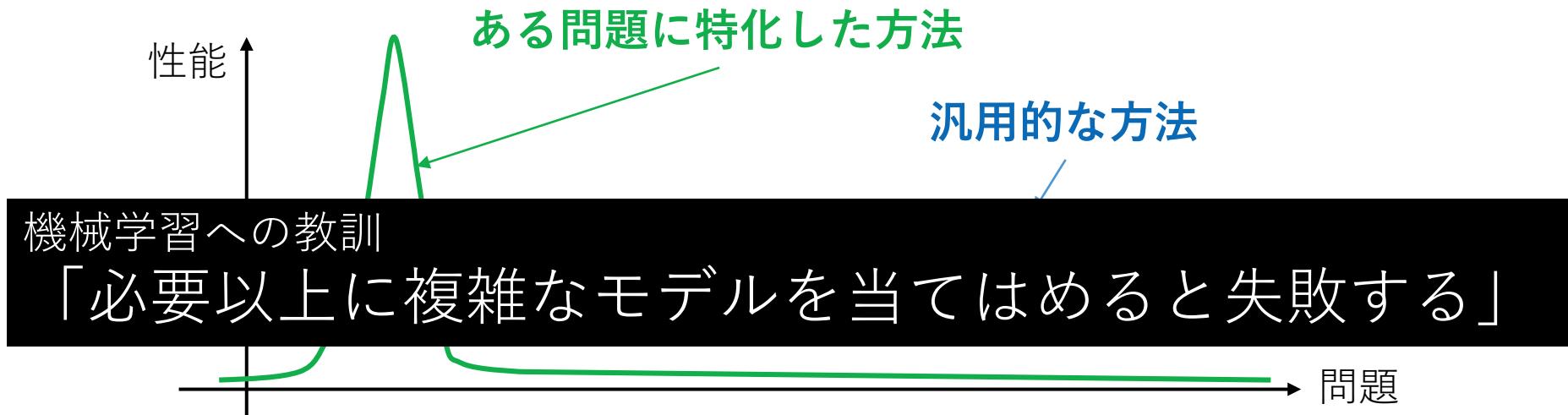
# 学習機の複雑さと学習能力

- オッカムの剃刀

「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」とする指針

- No free lunch theorem

「あらゆる問題で性能の良い汎用的学習機は実現不可能であり、ある問題に特殊化された手法に勝てない」

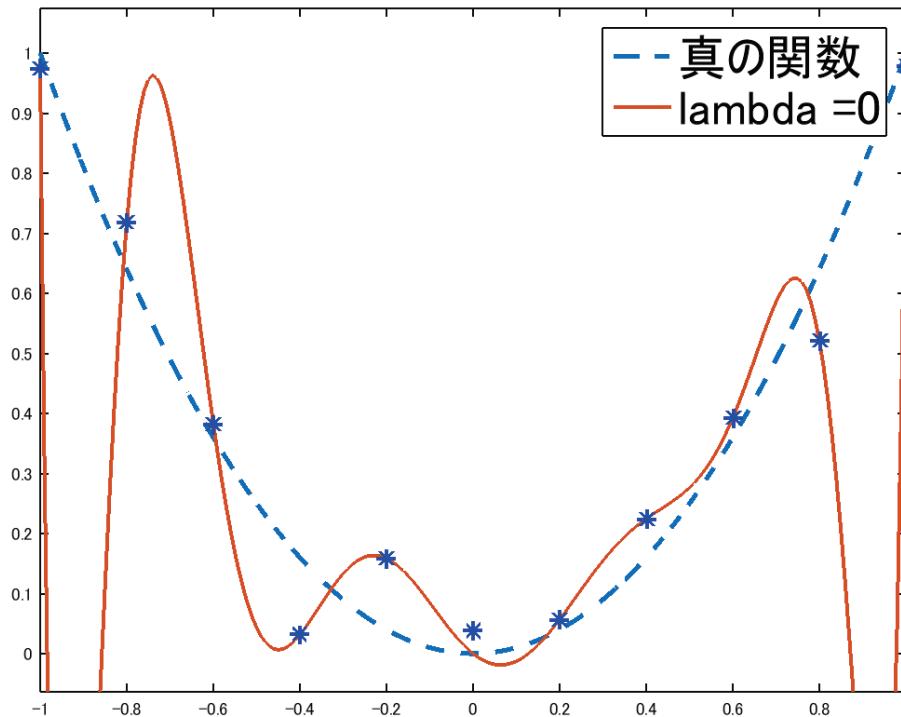


William of Ockham : 1285-1347. スコラ学の神学者, 哲学者.

No free lunch theorem: [D.H.Wolpert and W.G. Macready: 1995,1997][Y.C. Ho and D.L. Pepyne: 2002]

# 多項式回帰の過学習

$$\min_{\beta \in \mathbb{R}^{15}} \sum_{i=1}^n \{y_i - (\beta_1 x_i + \beta_2 x_i^2 + \cdots + \beta_{15} x_i^{15})\}^2$$



# 正則化学習法

正則化：データに合った単純なモデルを当てはめる  
→ 過学習を回避

正則化訓練誤差最小化

$$\min_{\beta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (y_i - x_i^\top \beta)^2 + \lambda \psi(\beta)$$

手元にあるデータへの当てはまり

正則化項

複雑さへの罰則

---

代表的な例：リッジ正則化 (L<sub>2</sub>ノルム)

$$\psi(\theta) = \|\theta\|_2^2 = \sum_{j=1}^d \theta_j^2$$

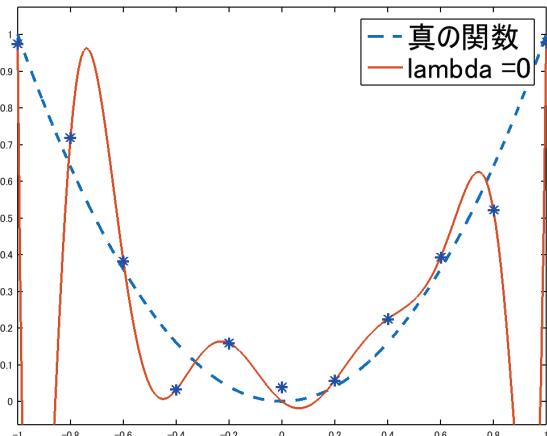
# 正則化の代表例

多項式回帰（15次多項式、リッジ回帰）

$$\min_{\beta \in \mathbb{R}^{15}} \sum_{i=1}^n \{y_i - (\beta_1 x_i + \beta_2 x_i^2 + \cdots + \beta_{15} x_i^{15})\}^2 + \lambda \|\beta\|_2^2$$

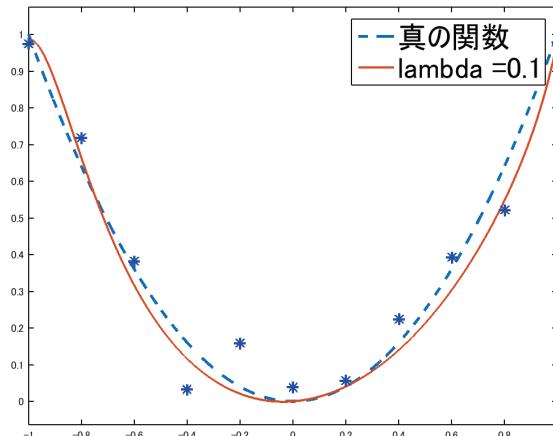
## リッジ正則化と言う

手元のデータには良くあてはまるが真の関数からは遠い



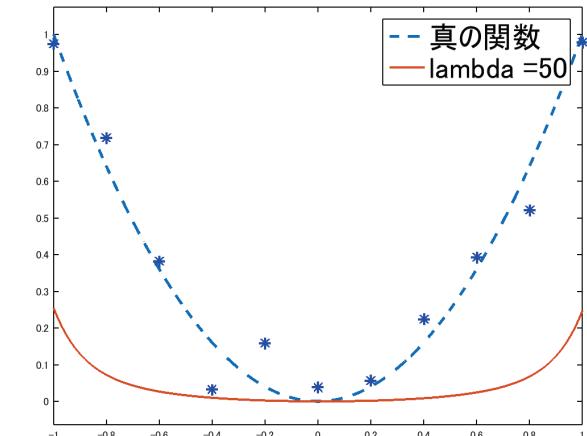
$\lambda = 0$

過学習



$\lambda = 0.1$

良い推定



$\lambda = 50$

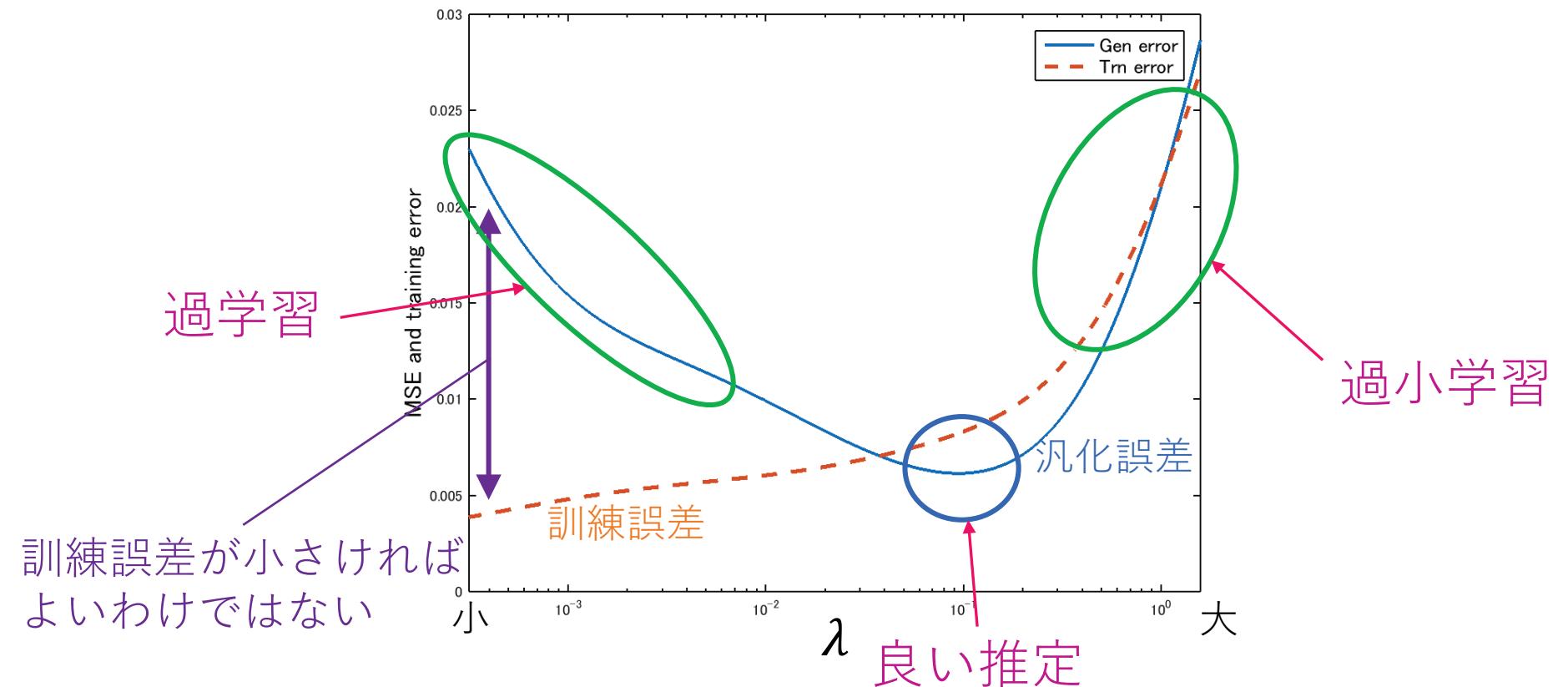
過小学習

正則化によってあまり複雑にならないよう制御がかかる

# 正則化の強さと汎化誤差の関係

多項式回帰 (15次多項式, リッジ回帰)

$$\min_{\beta \in \mathbb{R}^{15}} \sum_{i=1}^n \{y_i - (\beta_1 x_i + \beta_2 x_i^2 + \cdots + \beta_{15} x_i^{15})\}^2 + \lambda \|\beta\|_2^2$$



横軸：正則化パラメータ(log-スケール). 縦軸：汎化誤差（青）, 訓練誤差（赤）.

適切な  $\lambda$  を選ぶ方法→交差検証法, Mallows' Cp

# バイアスとバリアンスの分解

線形モデル (ノイズは平均0, 分散  $\sigma^2$ ) :

$$Y = X\beta^* + \epsilon$$

任意の推定量に対して以下の分解が成り立つ:

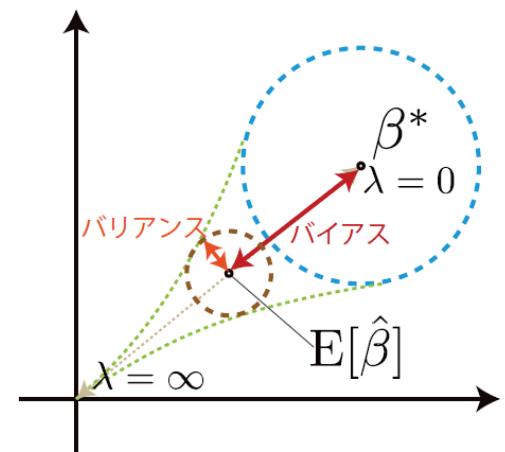
$$E[\|\hat{\beta} - \beta^*\|^2] = \underbrace{E[\|E(\hat{\beta}) - \beta^*\|^2]}_{\text{バイアス項}} + \underbrace{E[\|\hat{\beta} - E(\hat{\beta})\|^2]}_{\text{バリアンス項}}$$

リッジ正則化の場合 :  $\hat{\beta}_{(\lambda)} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ \sum_{i=1}^n (y_i - x_i^\top \beta)^2 + \lambda \|\beta\|^2 \right\}$

$$= \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \{ \|Y - X\beta\|^2 + \lambda \|\beta\|^2 \}$$

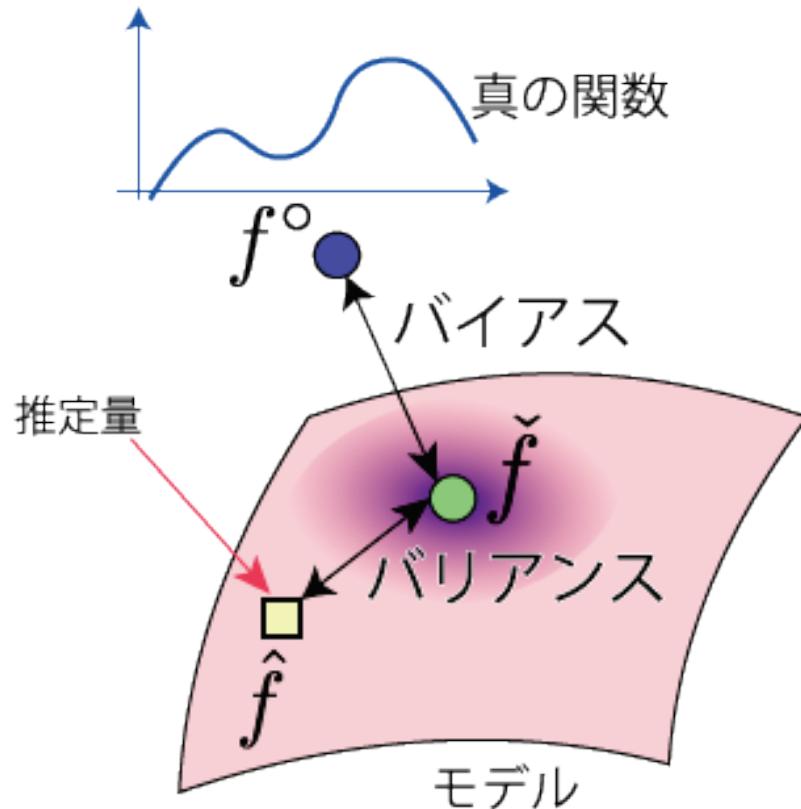
正則化パラメータとバイアス-バリアンス

	バイアス	バリアンス
$\lambda = 0$	0	$\sigma^2 \operatorname{Tr}[(X^\top X)^{-1}]$
$\lambda = \infty$	$\ \beta^*\ ^2$	0



※両方を同時に小さくすることはできない。

# バイアスとバリアンスの分解



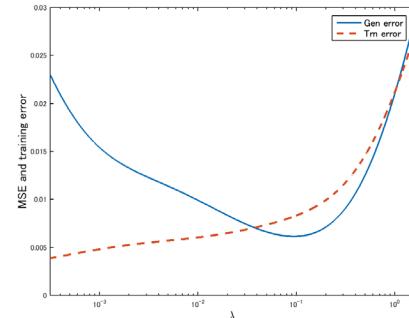
## バイアスとバリアンスのトレードオフ

- モデルが大きい：バイアス小，バリアンス大
- モデルが小さい：バイアス大，バリアンス小

サンプルサイズに合わせて適切なモデルを選択する必要がある。

# Mallows' CP規準

$$\begin{aligned}\hat{\beta}_{(\lambda)} &= \operatorname{argmin}_{\beta \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (y_i - x_i^\top \beta)^2 + \lambda \|\beta\|^2 \right\} \\ &= \operatorname{argmin}_{\beta \in \mathbb{R}^p} \{ \|Y - X\beta\|^2 + \lambda \|\beta\|^2 \}\end{aligned}$$



Mallows' CP規準

$$\hat{L}(\lambda) := \sum_{i=1}^n (y_i - x_i^\top \hat{\beta}_{(\lambda)})^2 + 2\hat{\sigma}^2 \operatorname{Tr}[X(X^\top X + \lambda I)^{-1} X^\top]$$

訓練誤差 補正項

$\hat{L}(\lambda)$  を最小にする  $\lambda$  を選択.

- Mallows' CP 規準  $\hat{L}(\lambda)$  は予測誤差  $E_{x,y}[(y - x^\top \hat{\beta})^2]$  の推定量.
- $\hat{\sigma}^2$  としては最小二乗推定量  $\hat{\beta}_{\text{LS}}$  を用いて  $\hat{\sigma}^2 = \|Y - X\hat{\beta}_{\text{LS}}\|^2/n$  を用いることが多い.

※  $\lambda = 0$  の時, AIC と一致する.

# クロスバリデーション (Cross-Validation)<sup>61</sup>

## 適切なハイパーパラメータを選ぶ方法

- 観測データへの当てはまりではなく予測誤差を最小化.
- 観測データへの当てはまりを最良にするのは  $\lambda = 0$ .

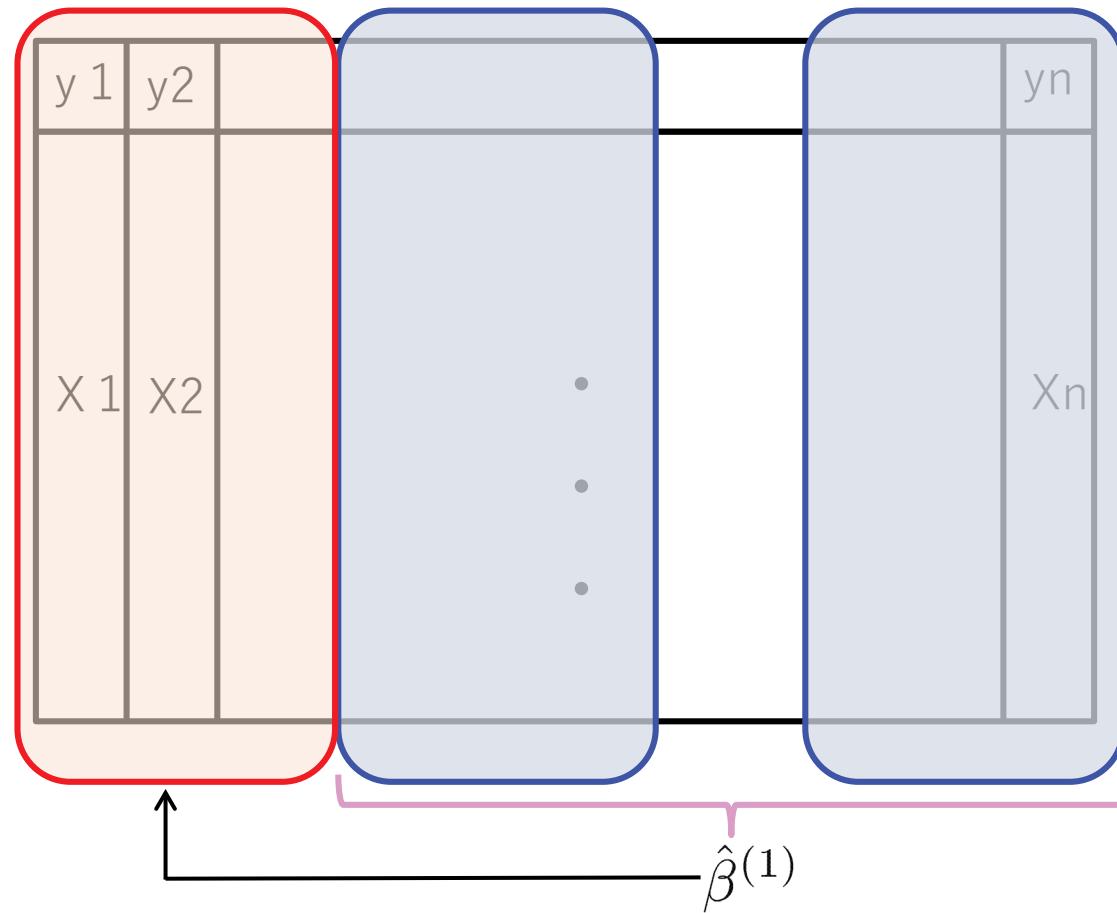
とにかくあらゆる問題に適用可能  
「とりあえずクロスバリデーション」

### **k-fold** クロスバリデーション

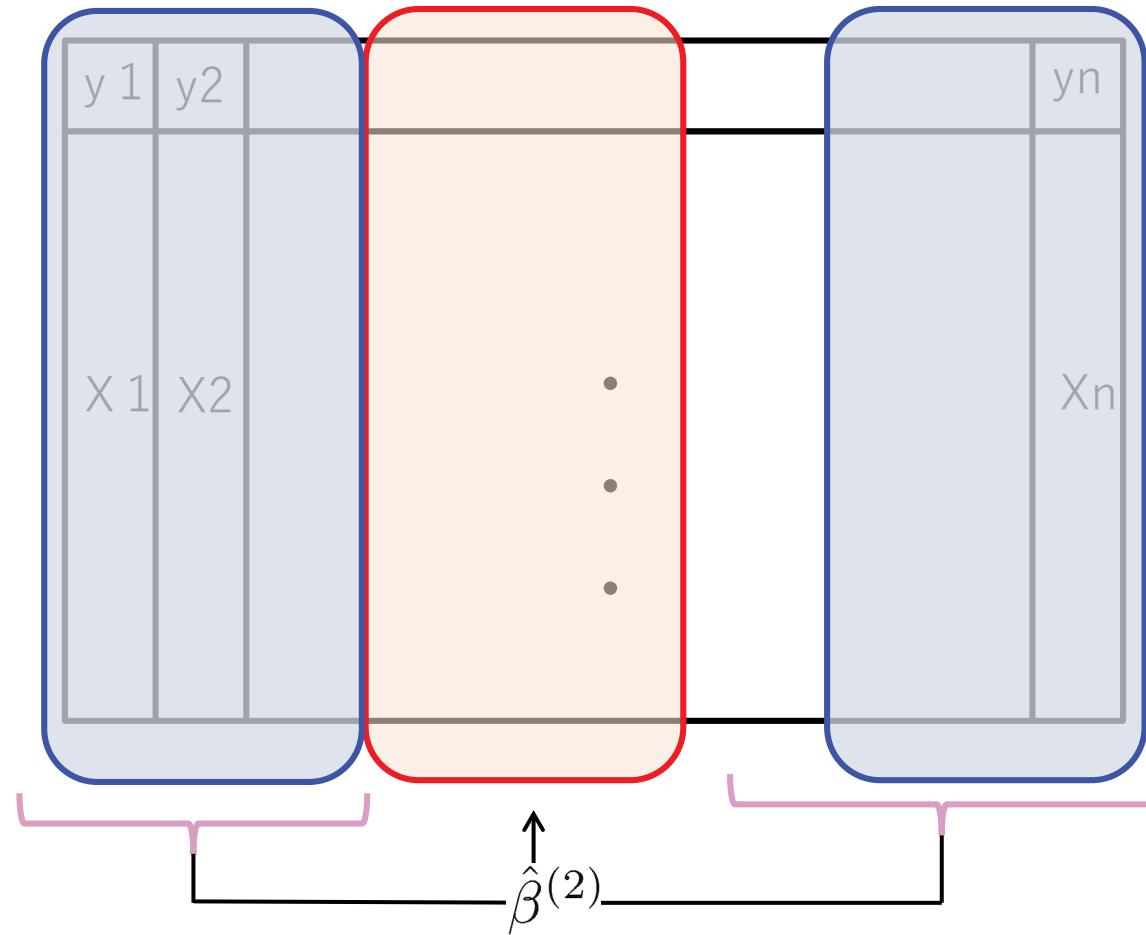
1. まずデータを  $k$  個に分割する.
2. 分割したデータの一つをテストデータとしてとっておき、残りのデータで推定.
3. テストデータでの予測誤差を計算.
4. 手順 2-3 を  $k$  個のテストデータの取り方について繰り返す.
5.  $k$  回繰り返しの予測誤差の平均を取る=CV スコア.

CV スコアを最小にする  $\lambda$  を選べば良い.

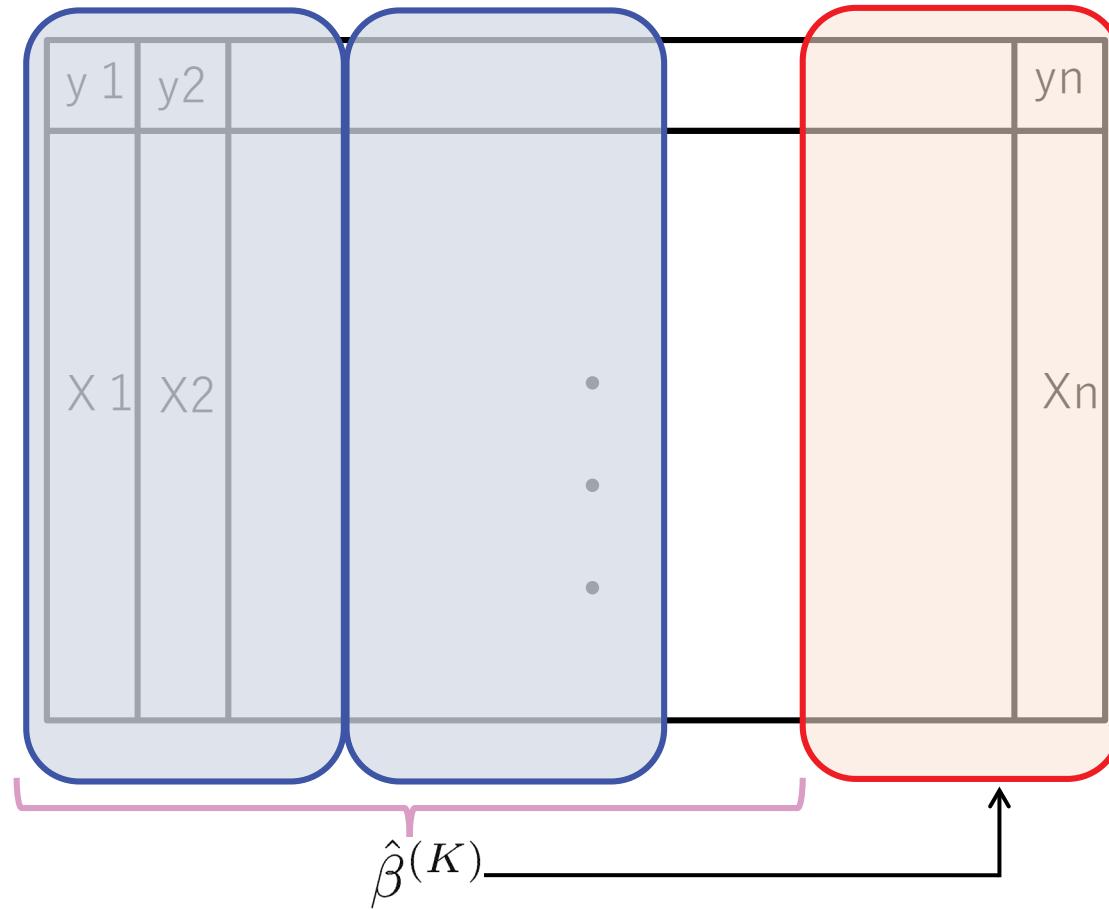
特に  $k = n$  (サンプルサイズ) の時、Leave-One-Out-CV (LOOCV) と呼ぶ.



$$\frac{1}{|I_1|} \sum_{i \in I_1} \ell(y_i, \hat{\beta}^{(1)\top} x_i)$$



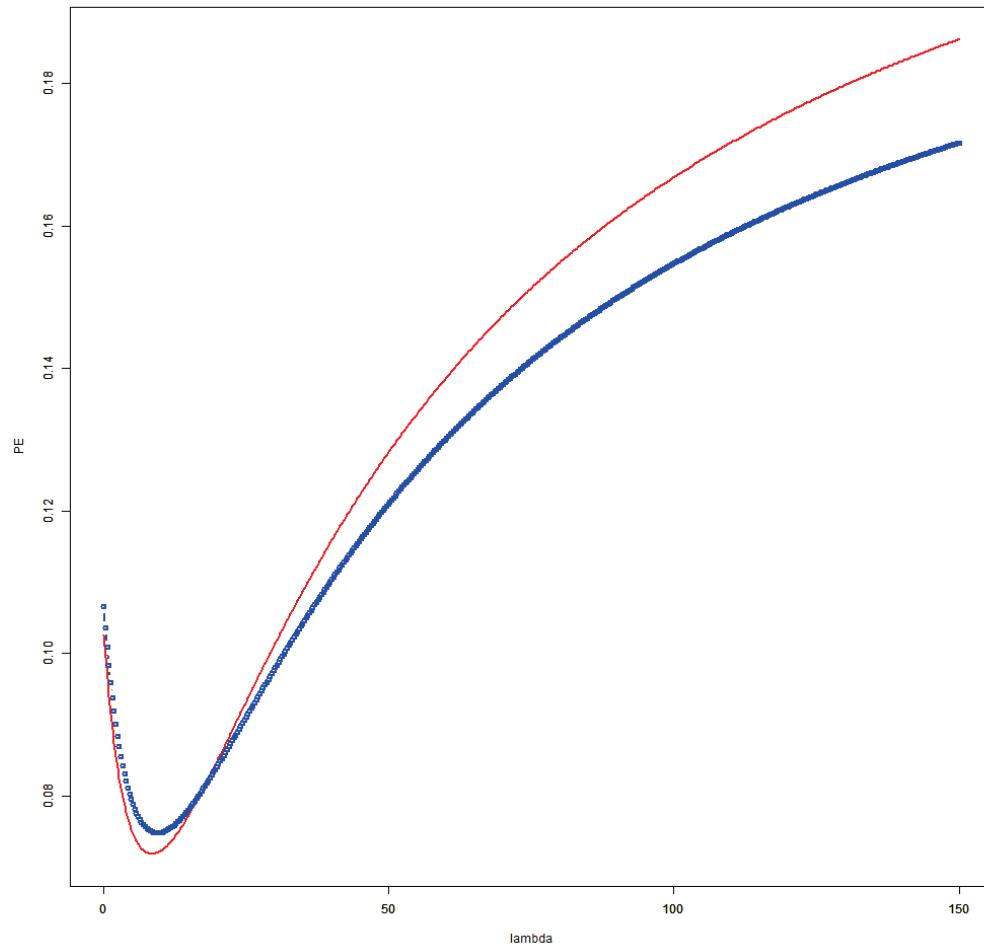
$$\frac{1}{|I_2|} \sum_{i \in I_2} \ell(y_i, \hat{\beta}^{(2)\top} x_i)$$



$$\frac{1}{|I_K|} \sum_{i \in I_K} \ell(y_i, \hat{\beta}^{(K)\top} x_i)$$

# 実例

$n = 100, d = 10$  のリッジ回帰 (ガウスマルコフモデル+二乗ノルム正則化)



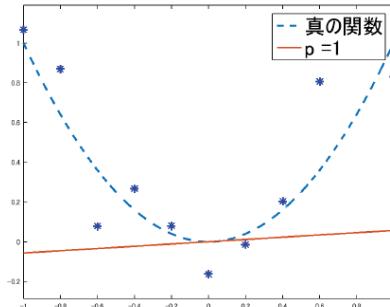
# 特徴選択

- 重回帰分析では説明変数を追加するごとに残差は小さくなつてゆく。
- 余分な説明変数を使うと過学習を起こしてしまう。
  - 観測済みデータによく当てはまつても、未観測のデータへの当てはまりが悪くなることがある。
  - サンプルサイズに比して複雑なモデルは使うべきではない。

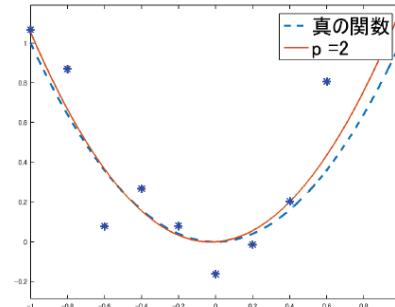
中古マンション価格の予測：床面積，築年数，駅からの距離，建ぺい率，…

例：多項式回帰

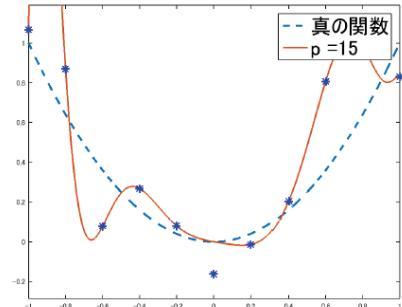
$$\min_{\beta \in \mathbb{R}^d} \sum_{i=1}^n \{y_i - (\beta_1 x_i + \beta_2 x_i^2 + \cdots + \beta_d x_i^d)\}^2$$



$d = 1$



$d = 2$



$d = 15$

# AICによる特徴選択

赤池情報量規準, AIC (Akaike Information Criterion)

予測精度が一番良いモデルを選択するための規準

$$\hat{\beta}^{(d)} = [\hat{\beta}_1, \dots, \hat{\beta}_d, 0, \dots, 0]^\top : d\text{変数のみを用いた最小二乗推定量}$$

$$AIC = \frac{\|Y - X\hat{\beta}^{(d)}\|^2}{\sigma^2} + 2d$$

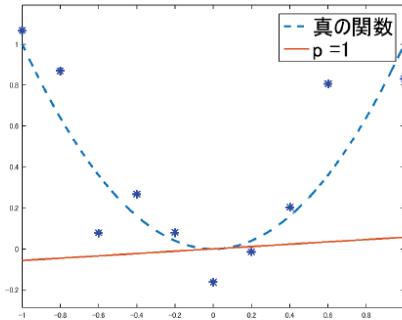
AICを最小化する変数の組を探す.

AIC = データへの当てはまりの良さ + モデルの複雑さ

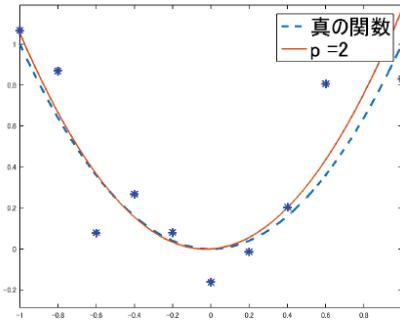
- $d$ 変数を用いた推定は最初の $d$ 変数である必要はない.
- $d$ を増やせばAICの第一項は減少し, 第二項は増大.
- AICの期待値は予測誤差になることが示せる.

# 例：多项式回帰

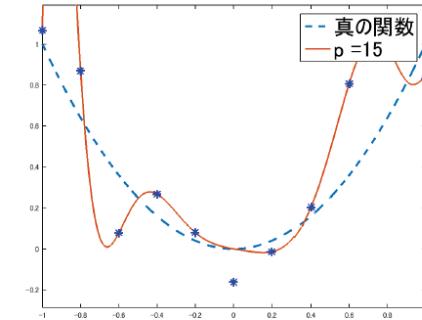
$$\min_{\beta \in \mathbb{R}^d} \sum_{i=1}^n \{y_i - (\beta_1 x_i + \beta_2 x_i^2 + \cdots + \beta_d x_i^d)\}^2$$



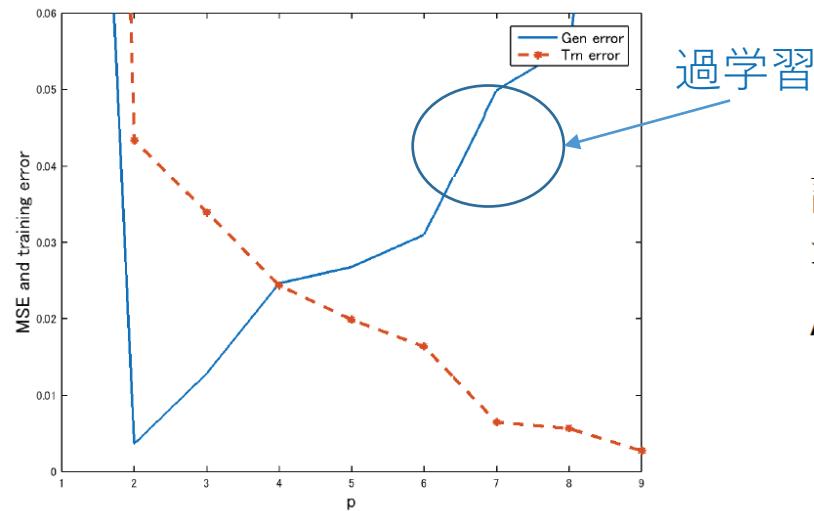
$$d = 1$$



$$d = 2$$



$$d = 15$$



訓練誤差は単調に減少するが、汎化誤差は途中で増大する。  
AICにより適切な次元が選ばれる。

汎化誤差と次元dの関係

# 中古マンション価格データの分析

69

国土交通省が公開している不動産取引価格情報から世田谷区の中古マンション取引価格データ（平成25年度第3四半期分）を取得。ここから一部を抜粋したデータで回帰分析をやってみる。

<http://www.land.mlit.go.jp/webland/download.html>

従属変数：価格

説明変数：1. 最寄駅からの距離（徒歩）

2. 延床面積

3. 建物の構造

4. 建ぺい率

5. 容積率

6. 建築年

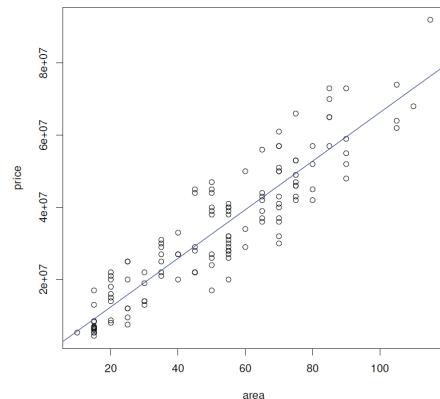
7. 最寄り駅に急行が止まるか（0-1変数で表現）

Rの関数 lm を使って分析。

# 回帰分析関数(lm)

## 最小二乗法の計算

```
sman.lm <- lm(price ~ area,data=sman) #回帰分析はこの一行でOK
plot(sman$area,sman$price, xlab="area",ylab="price") #結果をプロット
abline(sman.lm , lwd=1 , col="blue")
```



## AICによる特徴選択

```
sman.lmall <- lm(price ~ ., data=sman)
sman.lmAIC <- step(sman.lmall)
summary(sman.lmAIC)
```

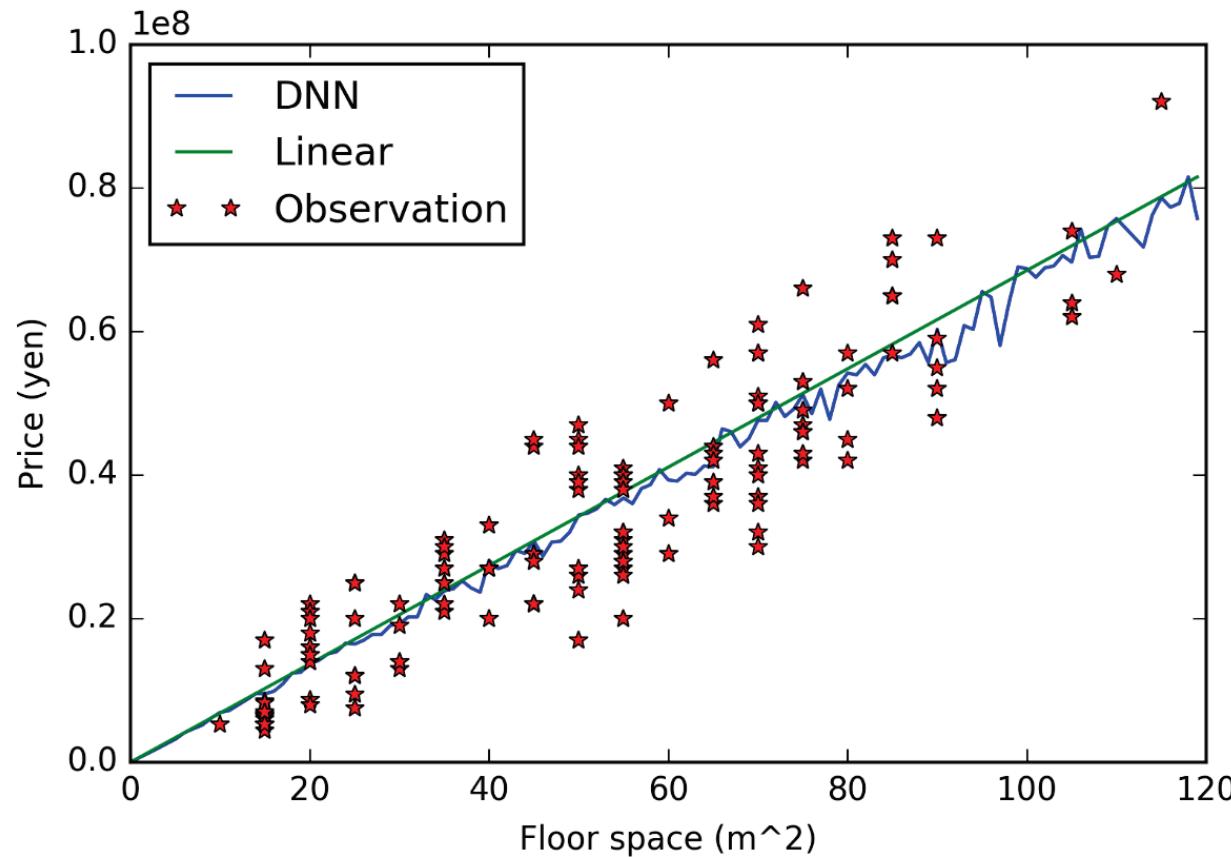
step() でAIC 最小のモデル(説明変数の組) を探索.

最寄駅からの距離 + 床面積 + 築年数  
の三変数モデルが採用された。

1. 最寄駅からの距離 (徒歩)
2. 延床面積
3. 建物の構造
4. 建ぺい率
5. 容積率
6. 築年数
7. 最寄り駅に急行が止まるか

# 線形モデル vs 深層学習

## 過学習の例



深層学習を使  
うには簡単&  
データが少な  
すぎる

マンションの価格推定

**DNN:** 中間層 2 層横幅100の深層NN, **Linear:** 単回帰モデル

汎化誤差（平均二乗誤差）: DNN:  $1.30 \times 10^{15}$ , Linear:  $6.26 \times 10^{13}$

一概に何でもかんでも深層学習が良いとは言えない

# これまでのまとめ

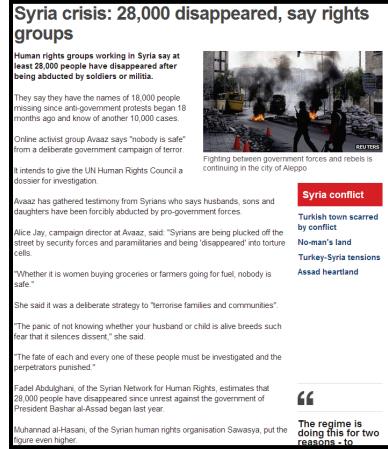
- 機械学習の歴史
- 機械学習の考え方
  - 複雑な規則をデータから学ぶ
- モデルと損失
  - 学習：期待損失最小化
- 過学習の問題
  - 複雑なモデルを当てはめれば良いわけではない.
  - 正則化
  - 変数選択

# 高次元スパース推定

# 高次元データ

インターネットや計測機器の発達により多様なデータが取得可能  
多くの場合で高次元

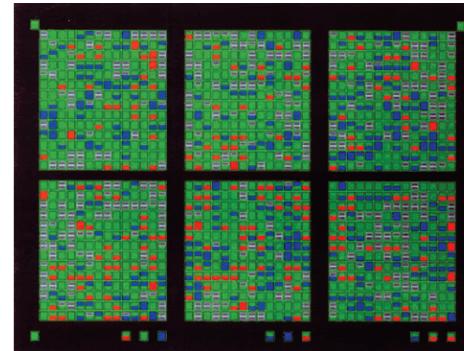
- 遺伝子データ
- テキストデータ
- マーケティングデータ
- 金融データ



Bag of words  
数百万次元

Syria	13
people	5
bomb	7
economy	1
immigrants	2
soccer	0
walk	1

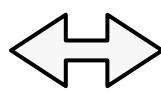
遺伝子発現量  
数万次元



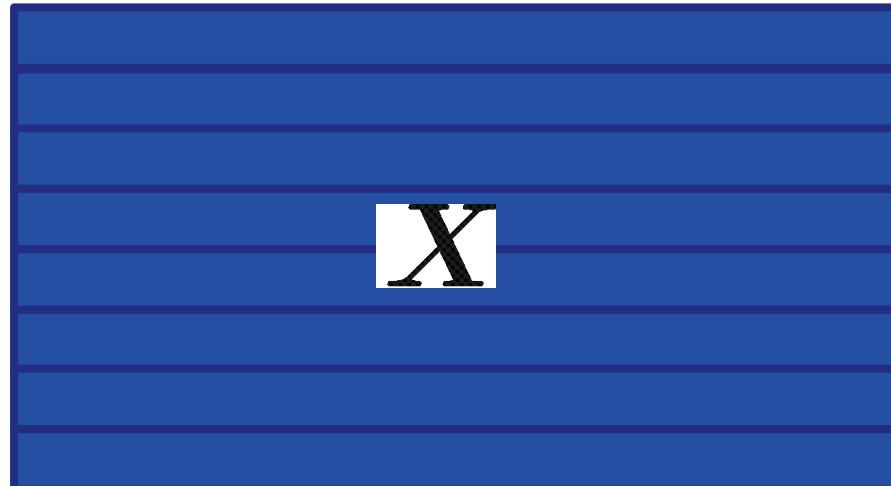
0.5
2.4
4.2
0.2
1.3
0.1
5.3

$y$  $x^\top$  $\beta$ 

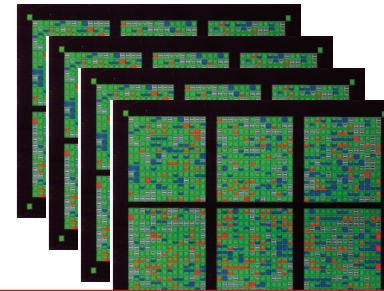
0.3
1.2
2.2
1.5
-0.5
-1.2
0.1
0.9



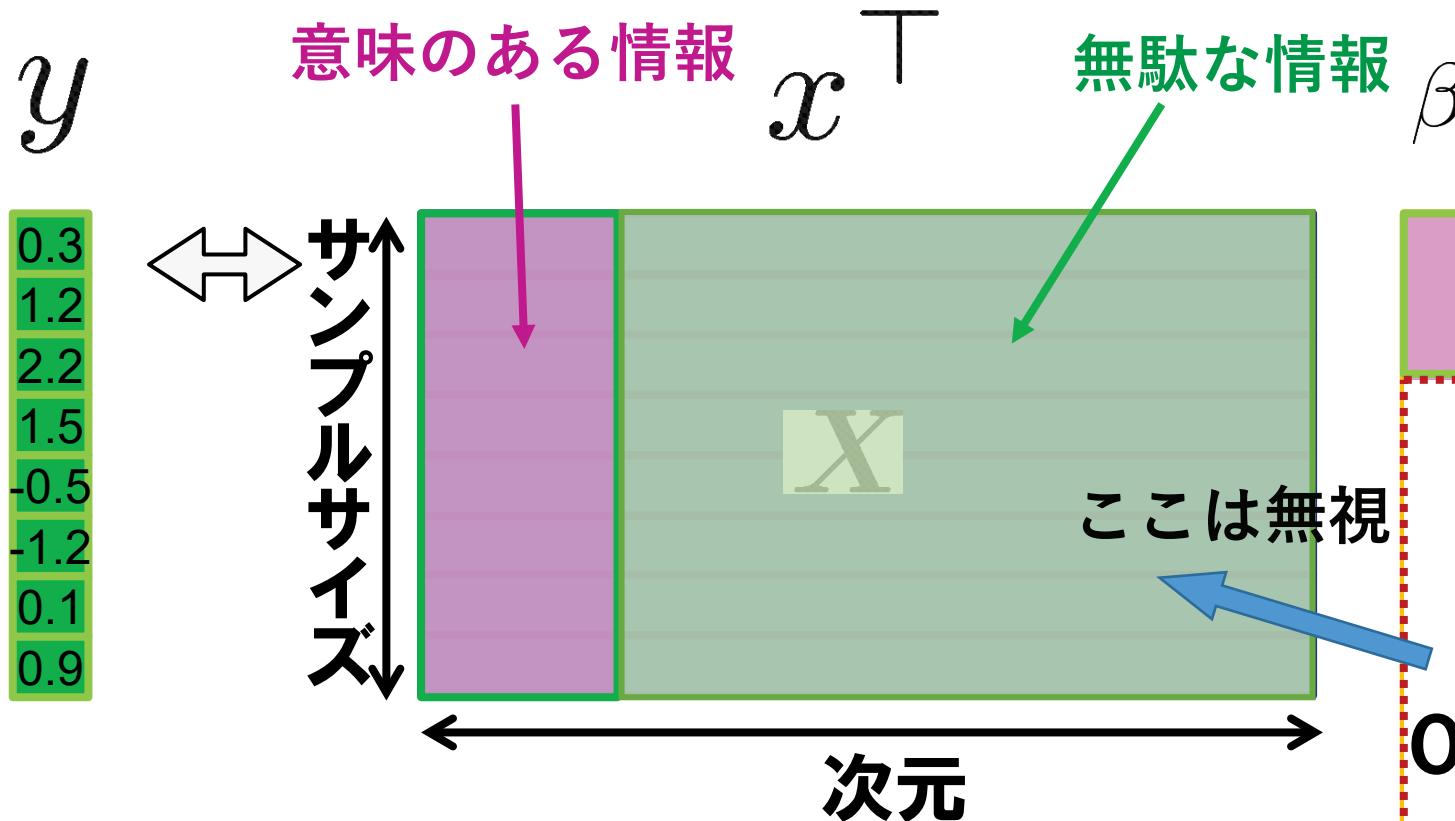
サンプルサイズ



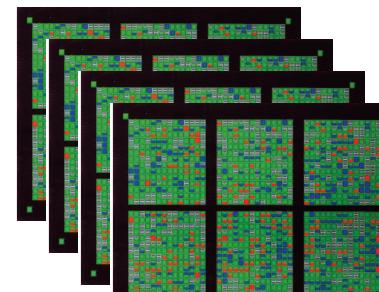
$\{(x_i, y_i)\}_{i=1}^n$ : サンプル



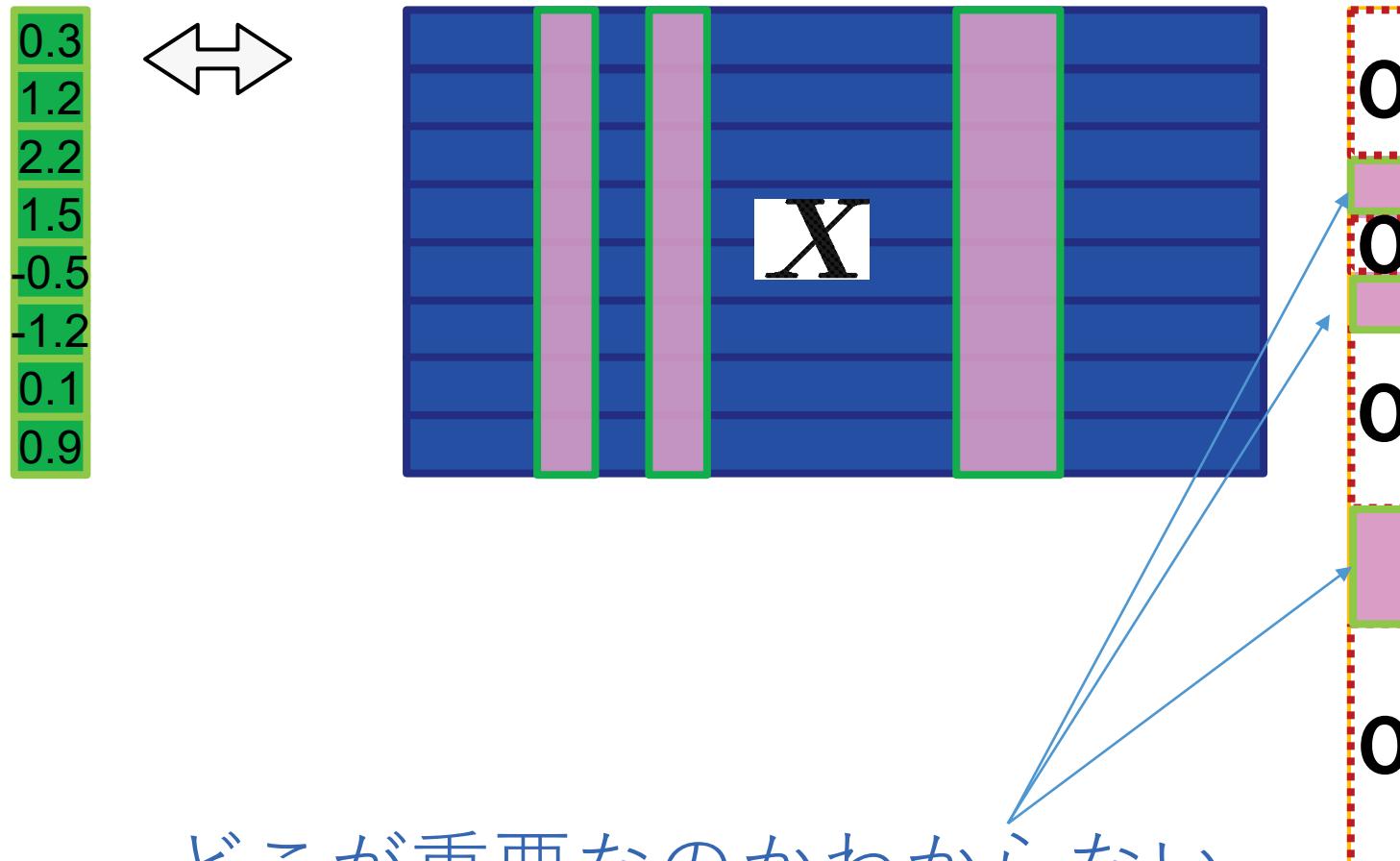
次元 > サンプルサイズ → 余分な情報を落としたい



$\{(x_i, y_i)\}_{i=1}^n$ : サンプル



次元 > サンプルサイズ → 余分な情報を落としたい  
スペースモデリング



どこが重要なのかわからない

→ 特徴選択：データから学習

予測に寄与する特徴量を特定できれば解釈性も上がる

# AICによる特徴選択（組み合わせ的方法）<sup>78</sup>

AIC: 赤池情報量規準 → 最尤推定量の予測誤差の不偏推定量

## AIC最小化

$$\hat{\beta}_{\text{AIC}} = \arg \min_{\beta \in \mathbb{R}^p} \|Y - X\beta\|^2 + 2\sigma^2 \|\beta\|_0$$

データへの  
当てはまり

次元に対する罰則  
(正則化)

ただし  $\|\beta\|_0 = \beta$  の非ゼロ要素の個数 :  $L_0$  ノルムと言う。

- 予測誤差を近似的に最小化
- 変数の組み合わせの数 :  $2^p$  個の候補 (膨大)
- NP困難

線形モデルを仮定

$$Y = X\beta^* + \xi$$

サンプルサイズ  $n$ , 次元  $p$

観測ノイズ : 分散  $\sigma^2$  の正規分布

# LASSOによる特徴選択（凸最適化）

**Lasso [ $L_1$ 正則化]** (R. Tibshirani (1996))

$$\hat{\beta}_{\text{Lasso}} = \arg \min_{\beta \in \mathbb{R}^p} \|Y - X\beta\|^2 + \lambda \|\beta\|_1$$

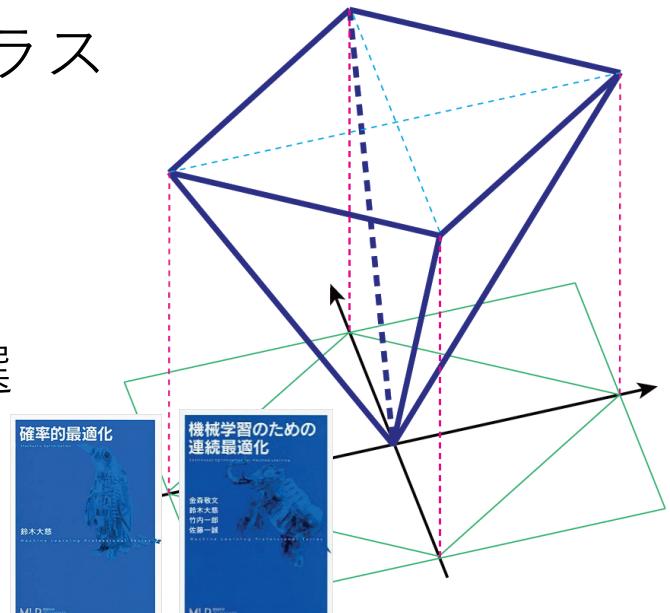
データへの  
当てはまり

次元に対する罰則  
(正則化)

ただし  $\|\beta\|_1 = \sum_{j=1}^p |\beta_j|$  :  $L_1$ ノルムと言う.

Lassoは**凸最適化**と呼ばれる問題のクラス

- 高速に解ける（近接勾配法等）
- $L_1$ ノルムは $L_0$ ノルムを最も良く近似する  
凸関数
- パラメータ $\lambda$ はクロスバリデーションで選  
べば良い。
- 理論が豊富。

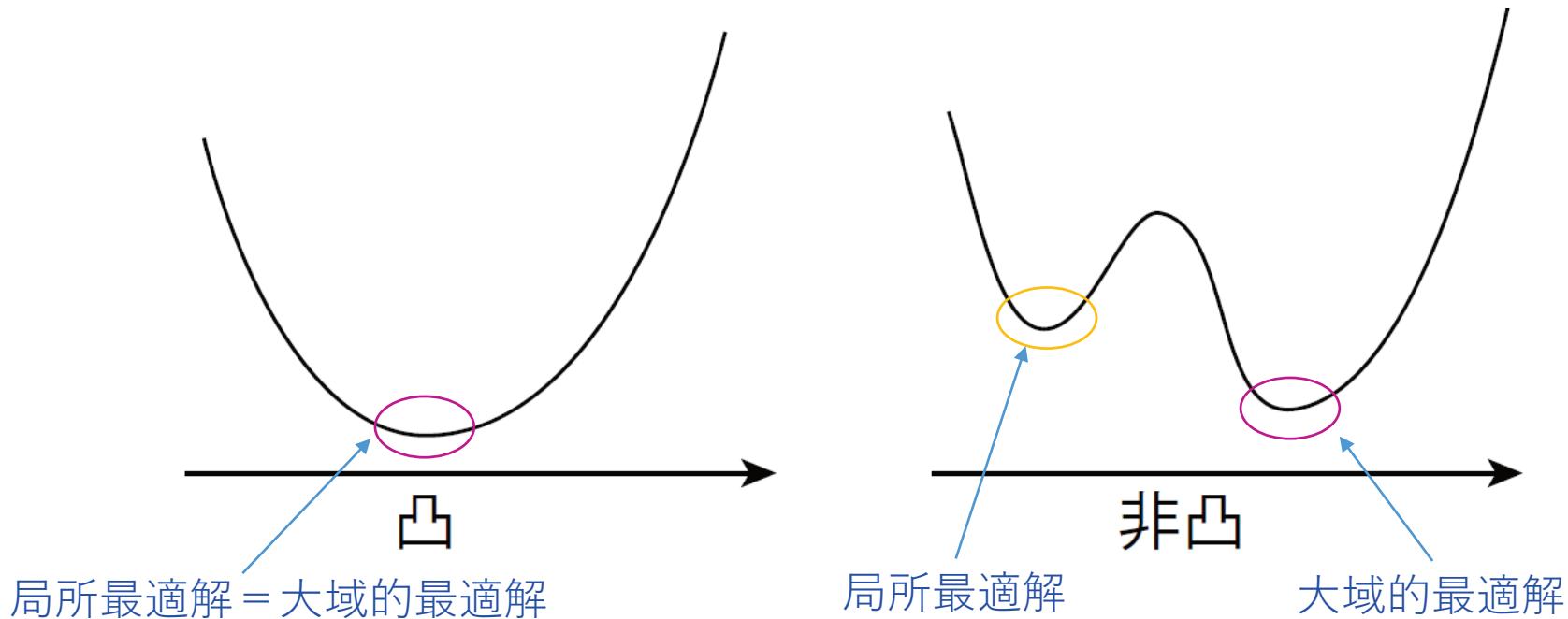


書籍：確率的最適化、機械学習のための連続最適化

# 凸関数

凸最適化 = 凸関数の最適化

$$\theta f(x) + (1 - \theta)f(y) \geq f(\theta x + (1 - \theta)y) \quad (\forall x, y \in \mathbb{R}^P, \theta \in [0, 1])$$



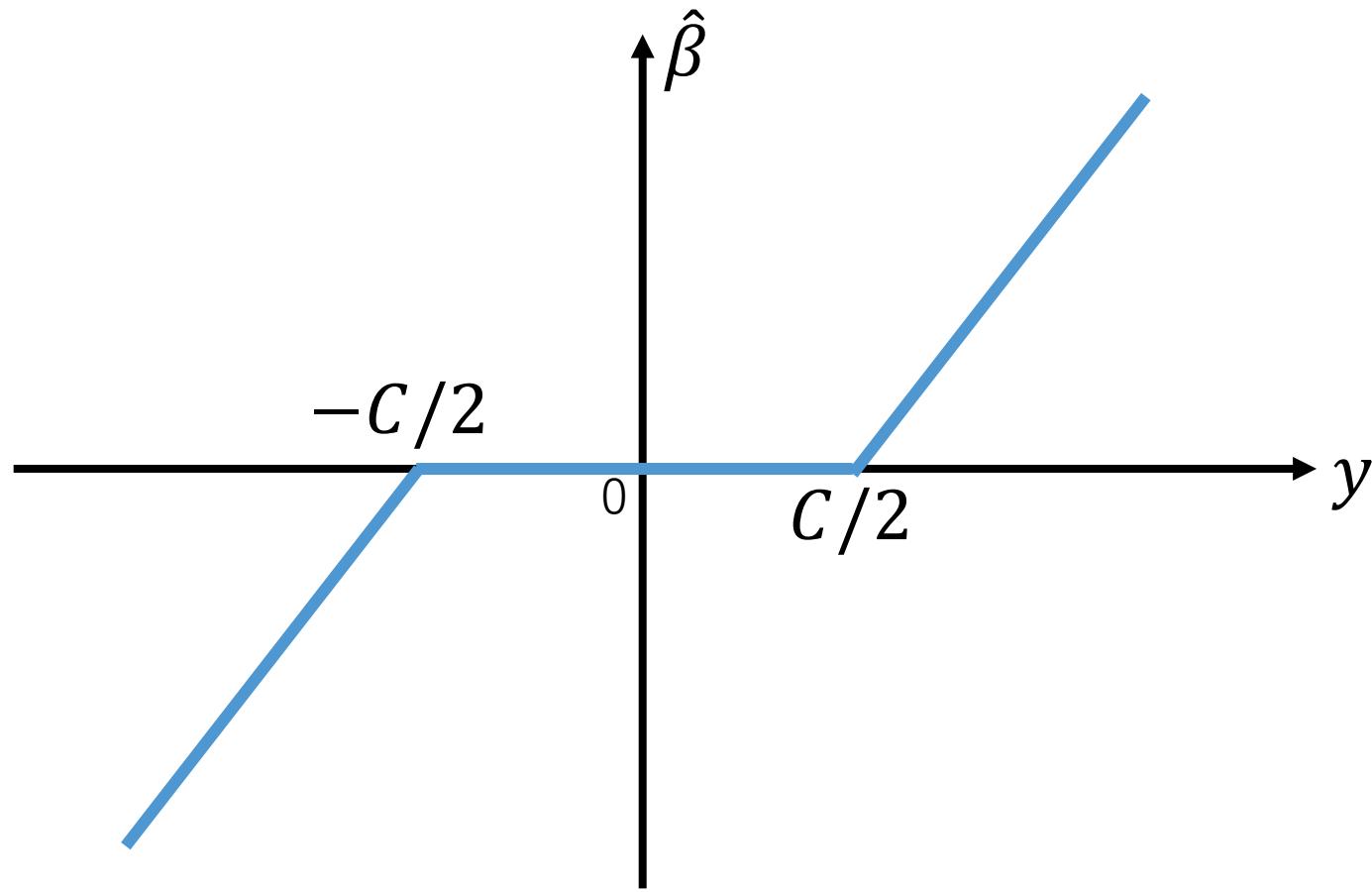
凸関数は局所最適解が大域的最適解

→ 効率的な最適化が可能な場合が多い

# 簡単な例

1次元の場合

$$\min_{\beta \in \mathbb{R}} (y - \beta)^2 + C|\beta|$$

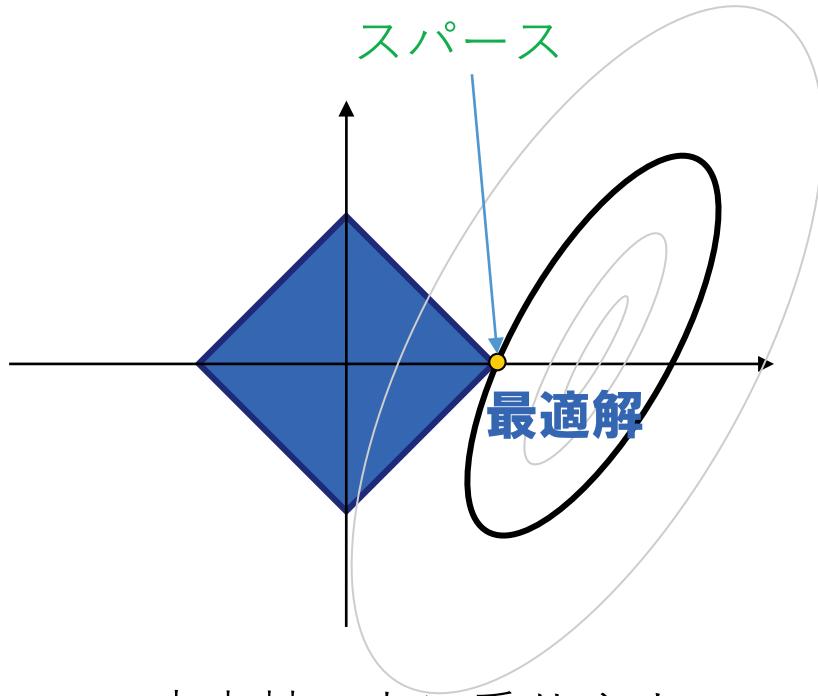


# Lassoのスパース性

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{n} \|X\beta - Y\|_2^2 + \lambda \sum_{j=1}^p |\beta_j|. \quad \leftrightarrow \quad \min_{\beta \in \mathbb{R}^p} \frac{1}{n} \|X\beta - Y\|_2^2 \text{ s.t. } \sum_{j=1}^p |\beta_j| \leq C$$

L1正則化

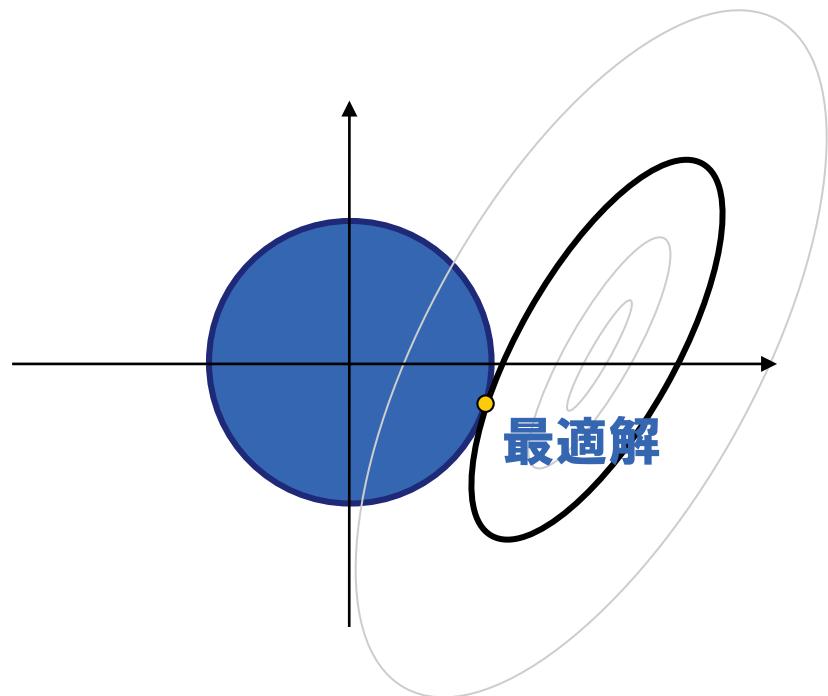
$$\|\beta\|_1 = |\beta_1| + \cdots + |\beta_p|$$



座表軸の上に乗りやすい

L2正則化（リッジ正則化）

$$\|\beta\|_2^2 = \beta_1^2 + \cdots + \beta_p^2$$



スパース推定によって予測に必要な変数が自動的に選ばれる

# スパース性の恩恵

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{n} \|X\beta - Y\|_2^2 + \lambda_n \sum_{j=1}^p |\beta_j|$$

$d$  = 真のベクトル  $\beta^*$  の非ゼロ要素の数 (予測に寄与する変数の数)

定理 (Lassoの収束レート (Bickel et al., 2009; Zhang, 2009))

ある条件のもと (制限等長性など) , ある定数  $C$  が存在して,

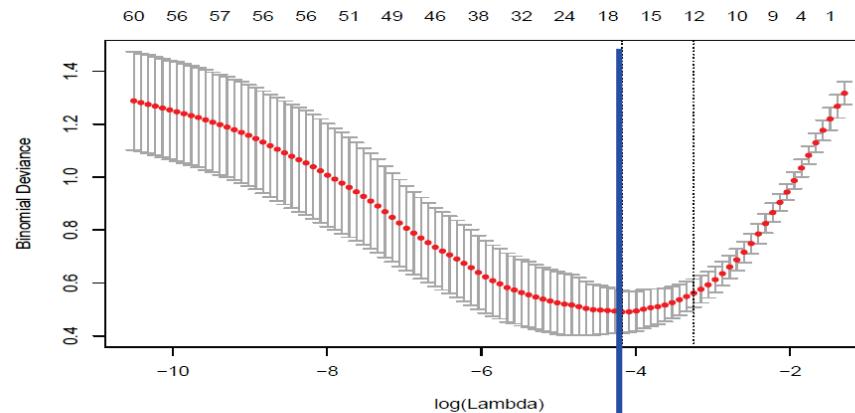
$$\|\hat{\beta} - \beta^*\|_2^2 \leq C \frac{d \log(p)}{n}$$

- 全体の次元  $p$  はたかだか  $O(\log(p))$  でしか影響しない!
  - 実質的次元  $d$  が支配的.
  - 高次元スパースな問題を精度よく解くことができる.
- 過学習を防止

推定誤差	$\frac{d \log(p)}{n} \ll \frac{p}{n}$	(最小二乗法)
		過学習してしまう

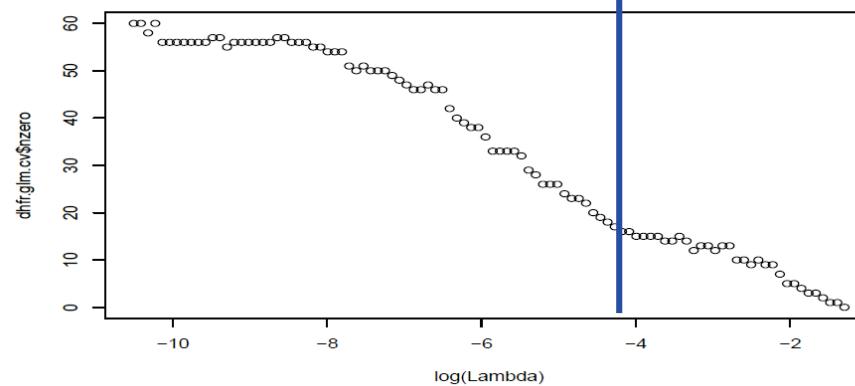
低次元性 (スパース性) をうまく利用できている.

## ジヒドロ葉酸レダクターゼデータにおける実験



CVスコア  
(予測精度)

ちょうどよい正則化

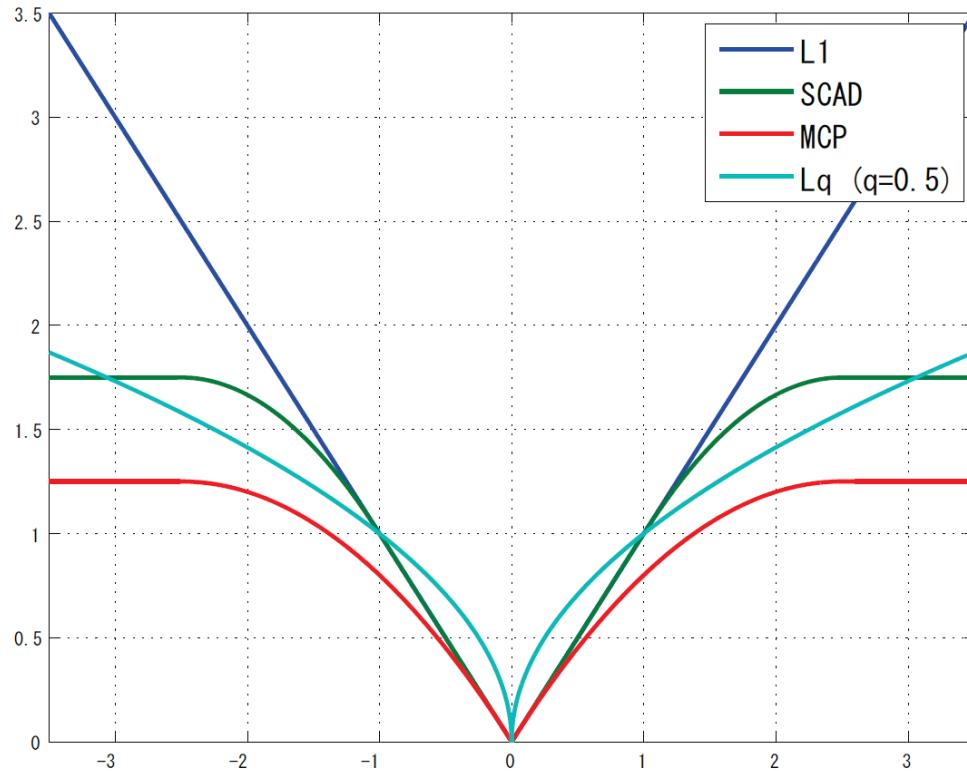


非ゼロ要素の個数

スパース性と汎化誤差

横軸：正則化パラメータ. 縦軸：(上段) CVスコア, (下段) 非ゼロ要素個数

# 非凸正則化



S C A D

$$\rho(|\beta|, \lambda) = \begin{cases} \lambda|\beta| & (|\beta| \leq \lambda) \\ \frac{-|\beta|^2 + 2a\lambda|\beta| - \lambda^2}{2(a-1)} & (\lambda < |\beta| \leq a\lambda) \\ \frac{(a+1)\lambda^2}{2} & (|\beta| \geq a\lambda) \end{cases}$$

M C P

$$\rho(|\beta|; \lambda) = \lambda \int_0^{|\beta|} (1 - x/(\gamma\lambda))_+ dx$$

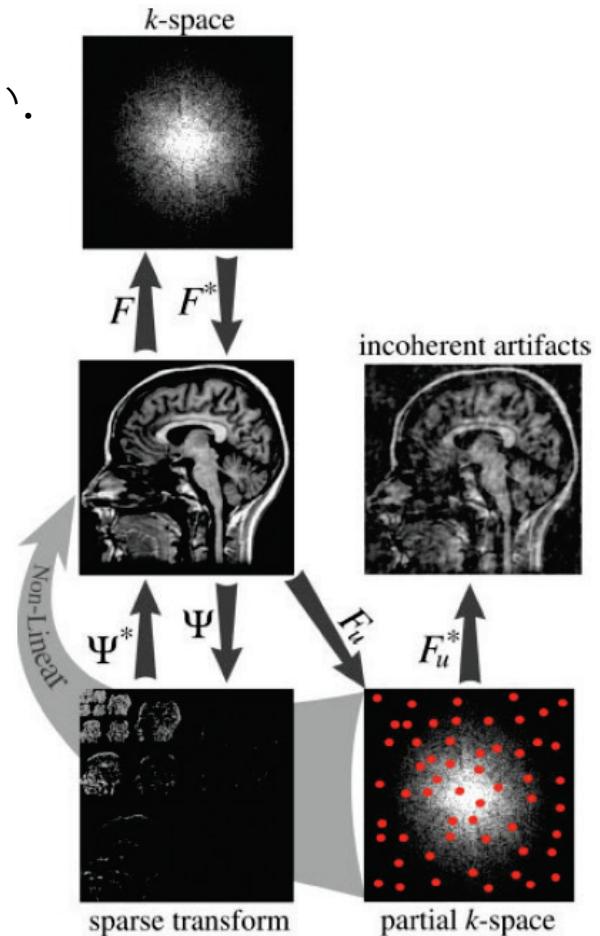
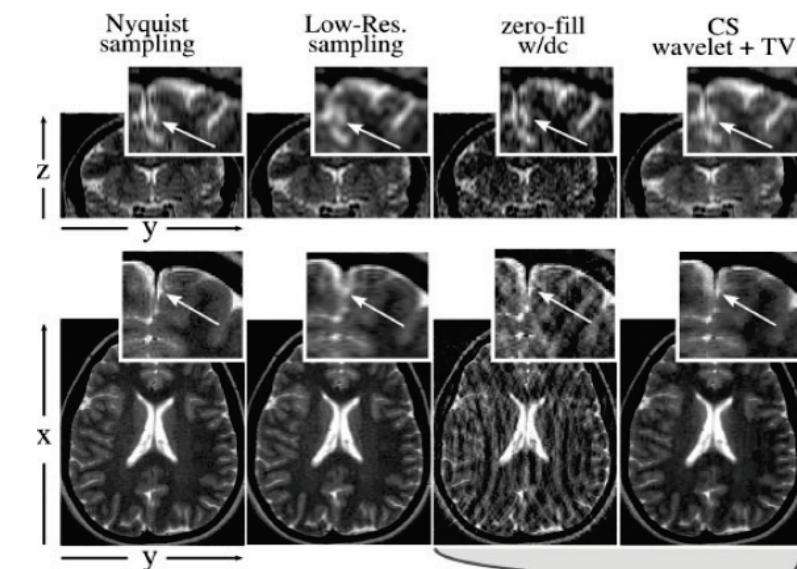
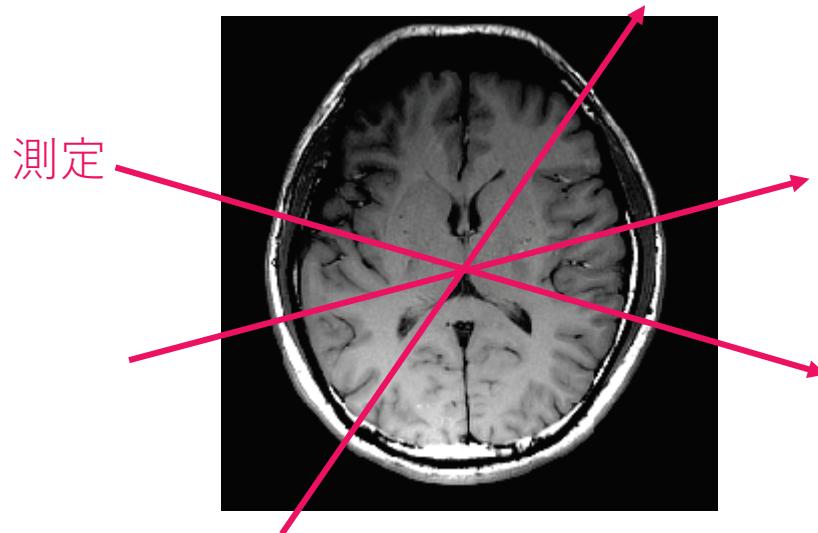
- SCAD (Smoothly Clipped Absolute Deviation) (Fan and Li, 2001)
- MCP (Minimax Concave Penalty) (Zhang, 2010)
- Lq 正則化( $q < 1$ ), Bridge 正則化(Frank and Friedman, 1993)

よりスパースな解. その代わり最適化は難しくなる.

ただし, 最近は局所最適解でも統計的性質は良いことが示されている.

# MRIへの応用

なるべく測定時間（観測回数）を減らしたい。



画像はwavelet基底に関してスパース  
→ 少数の観測（サンプル）でも大丈夫

[Lustig, Donoho and Pauly: Sparse MRI: The application of compressed sensing for rapid MR imaging, 2007]

# スペース共分散選択

$$x_k \sim N(0, \Sigma) \quad (\text{i.i.d.}, \Sigma \in \mathbb{R}^{p \times p}), \quad \widehat{\Sigma} = \frac{1}{n} \sum_{k=1}^n x_k x_k^\top$$

$$\hat{S} = \underset{S: \text{半正定対称}}{\arg \min} \left\{ -\log(\det(S)) + \text{Tr}[S \widehat{\Sigma}] + \lambda \sum_{i,j=1}^p |S_{i,j}| \right\}.$$

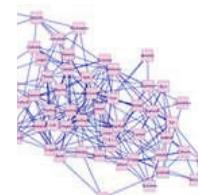
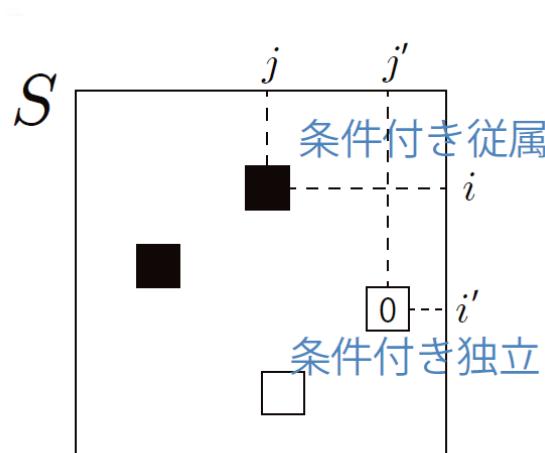
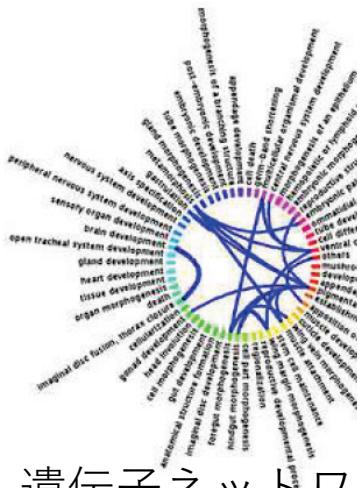
データへの当てはまり  
(正規分布の負対数尤度)

L1正則化

- の逆行列  $S$  を推定
- $S_{ij}=0 \Leftrightarrow$  「 $X_i$  と  $X_j$  が条件付き独立」
- $S_{ij}=0$  なら変数  $X_i$  と変数  $X_j$  は直接的に相互作用しないという意味

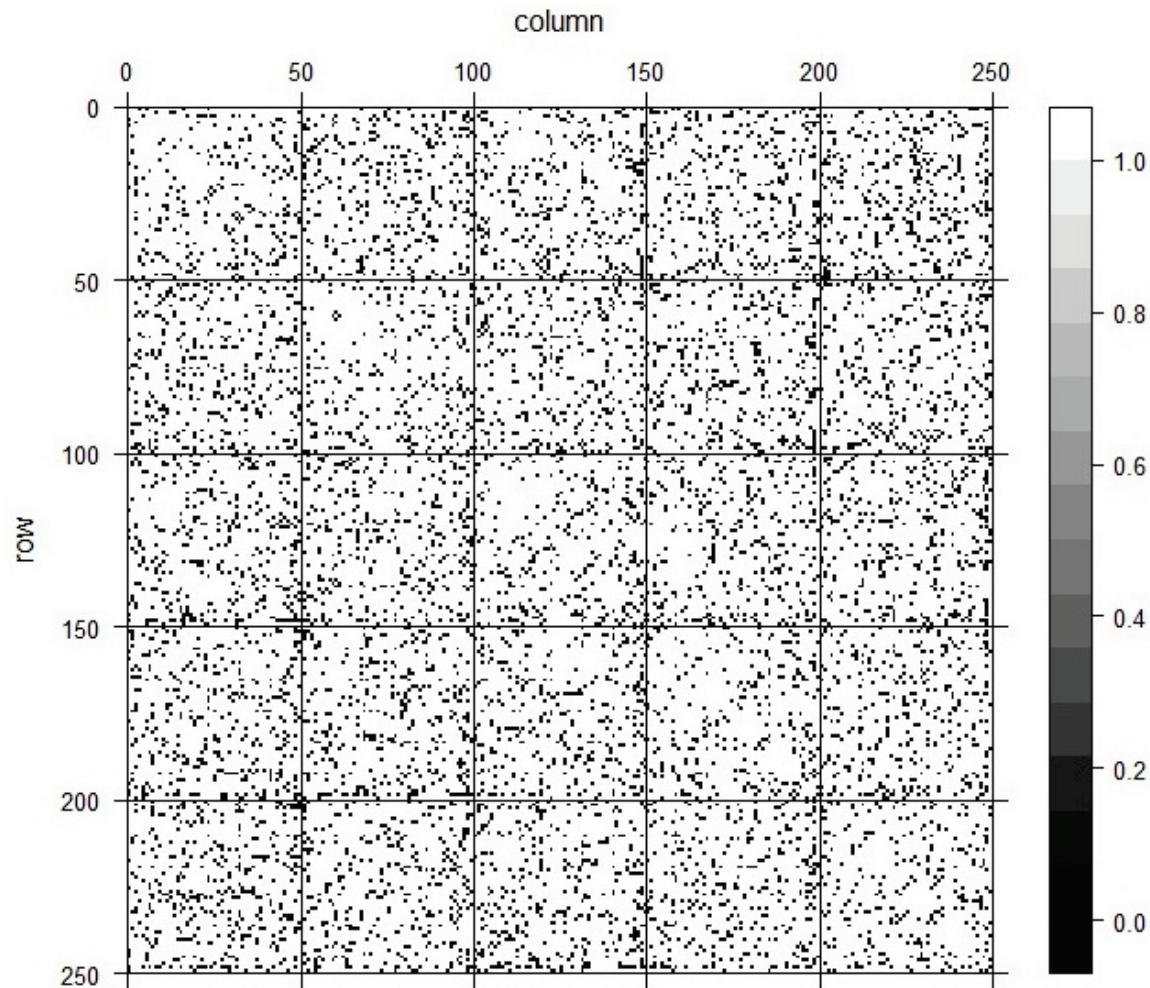
[Meinshausen and Bühlmann, 2006, Yuan and Lin, 2007, Banerjee et al., 2008]

グラフィカルモデルが 凸最適化で推定可能



員間の関係  
( $\Sigma$  から推定)

[Kolar+etal, 2010]



NASDAQ 銘柄からランダム抽出した50 銘柄.  
株価データを用いた分散共分散選択. 時間差も考慮.  
(2011 年1 月4 日から2014 年12 月31 日まで)  
(Lie Michael, Bachelor thesis)

# その他のスパース性

スパース正則化はL1正則化だけではない。  
他にも以下のようなより構造を持った正則化がある。

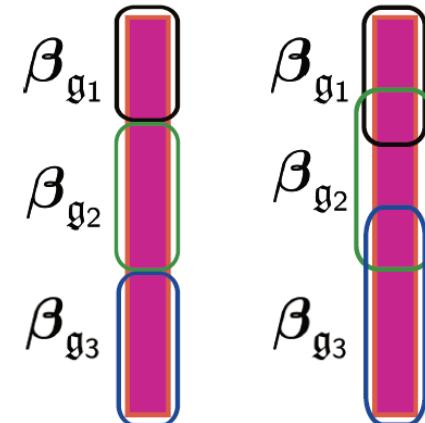
## 構造的正則化

- グループ正則化：変数のグループごと0にする。
- 一般化連結正則化
- トレスノルム正則化

# グループ正則化

$$\psi(\beta) = C \sum_{g \in \mathfrak{G}} \|\beta_g\|$$

- グループごとに正則化
- グループ全体が0になりやすい。



	SNP1	SNP2	SNP3	SNP4	SNP5	SNP6
Haplotype 1	CAGATCGCTGAATGAATCGCATTGT					
Haplotype 2	CAGATCGCTGAATGGATCCCATCAGT					
Haplotype 3	CGGATTGCTGCATGGATCCCATCAGT					
Haplotype 4	CGGATTGCTGCATGAATCGCATTGT					

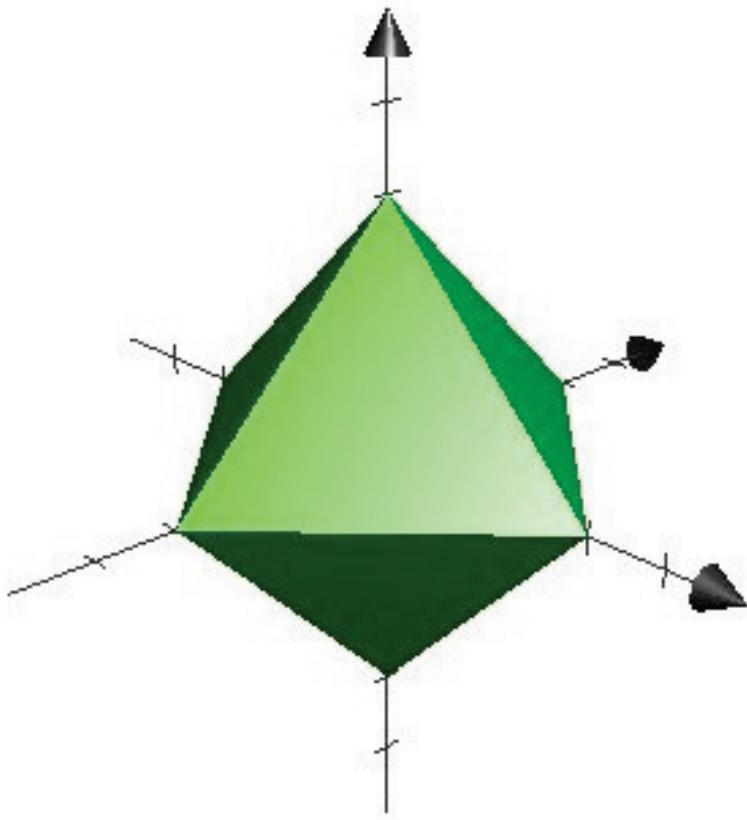
Group1                    Group2                    Group3

Genome Wide Association Study (GWAS)  
 [Balding '06, McCarthy et al. '08]

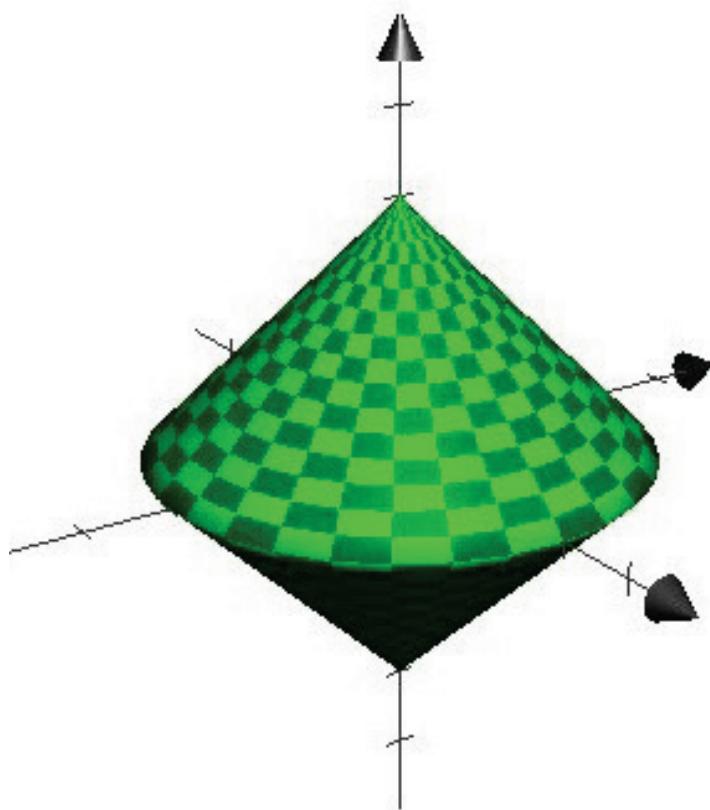
# グループ正則化の概形

91

Lasso



Group Lasso

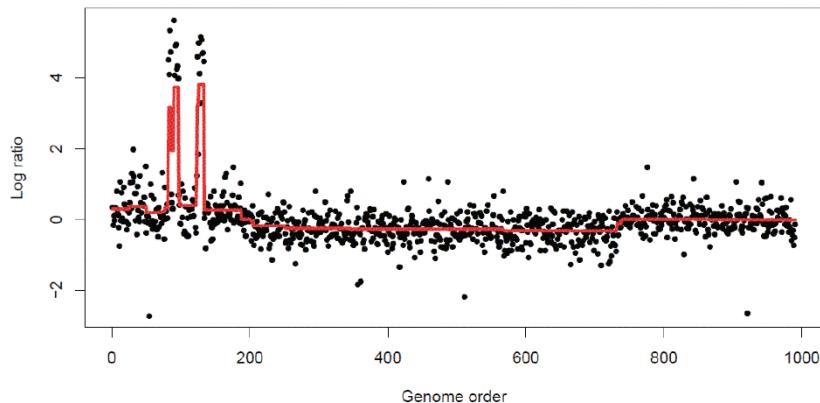


$$|\beta_1| + |\beta_2| + |\beta_3| \leq 1$$

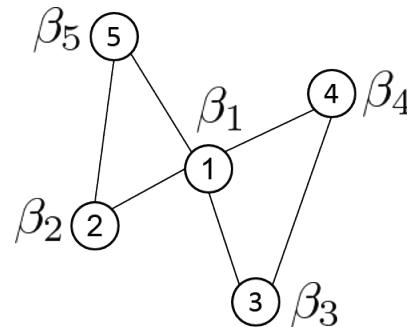
$$\sqrt{\beta_1^2 + \beta_2^2} + |\beta_3| \leq 1$$

# 一般化連結正則化 (Fused Lasso)

$$\psi(\beta) = \sum_{(i,j) \in E} |\beta_i - \beta_j|$$



Fused lasso による遺伝子データ解析  
[Tibshirani and Taylor '11]



TVデノイジング  
(パッチを使わないデノイジング)  
[Chambolle '04, Mairal et al., 2009]

背景切り出し [Mairal et al.: 2011]

テスト画像



L1正則化



L1/L2グループ正則化一般化連結正則化



# 低ランク行列補完

ベクトルから**行列**の学習へ

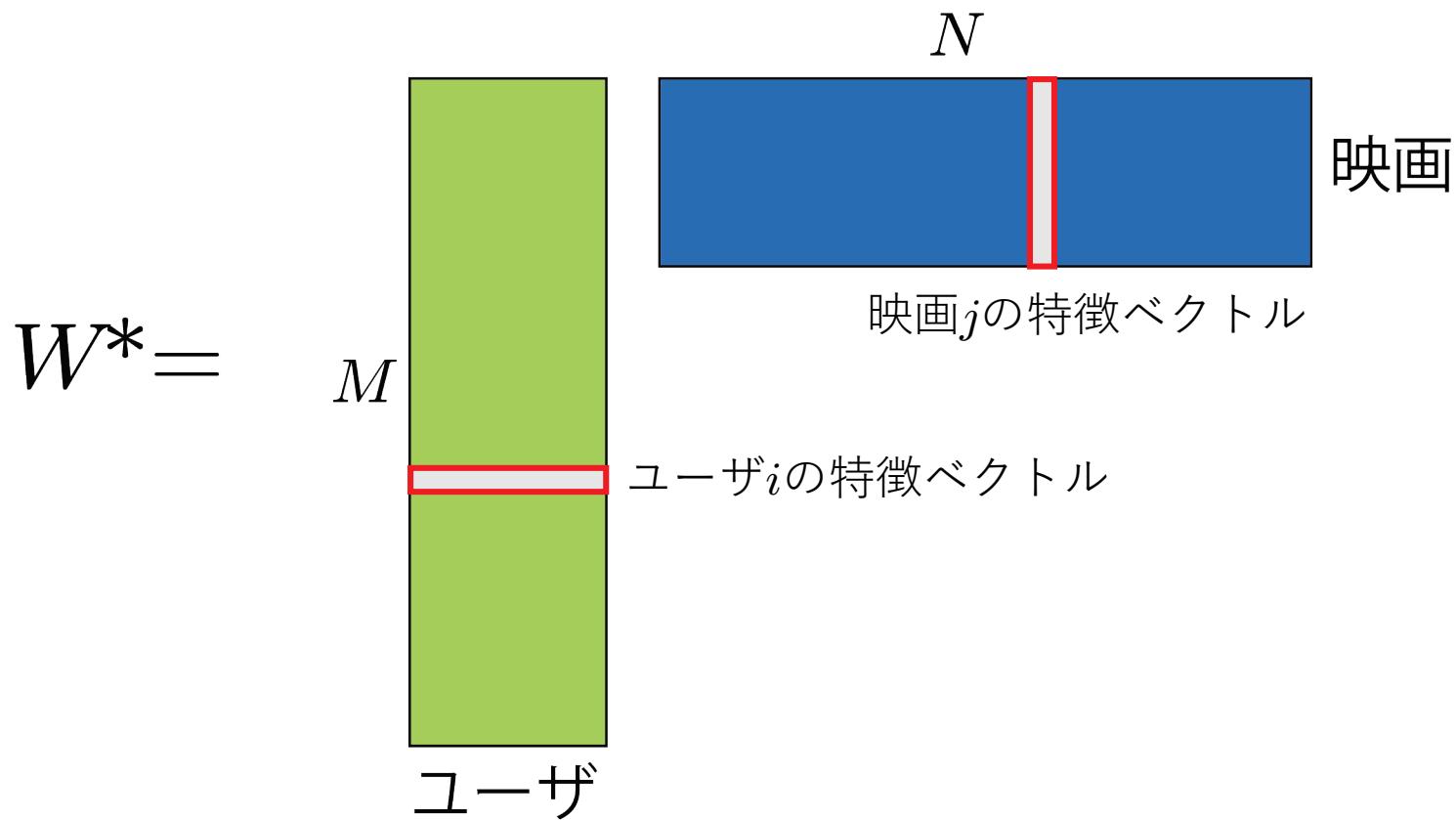
- 推薦システム

	映画 A	映画 B	映画 C	…	映画 X
ユーザ 1	4	8	4	…	2
ユーザ 2	2	4	2	…	1
ユーザ 3	2	4	2	…	1
:					

ランク 1 と仮定

各ユーザーが各映画をどれだけ好むかという部分的情報がある。  
 → 残りの部分 (\*の部分) を埋めたい。  
低ランク行列補完で可能。

e.g., Netflix prize (100万ドルの賞金, 48万ユーザ×1万8千映画)

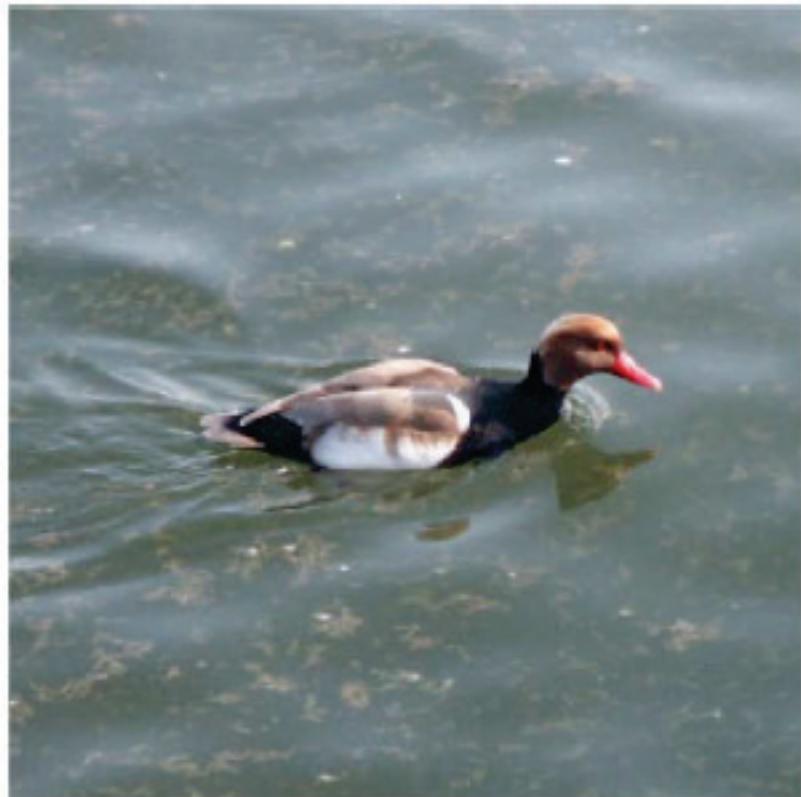


低ランク行列の学習は「ユーザ」と「映画」の**低次元表現**を学習することに他ならない。  
→交互最適化法やトレスノルム正則化法で学習可能

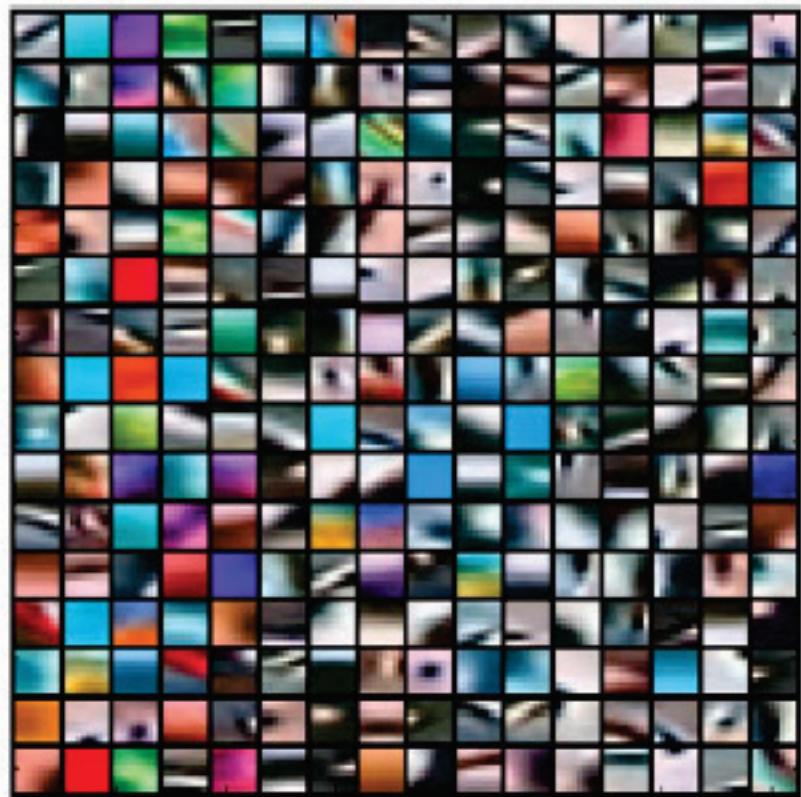
推定誤差  $O\left(\frac{r(M+N)}{n}\right) \ll O\left(\frac{MN}{n}\right)$  (低ランク性を利用しない最小二乗法)

$r$ : ランク

# スペース表現, 辞書学習



(a)



(b)

Mairal, Elad and Sapiro: Sparse Representation for Color Image Restoration.  
*IEEE Transactions on Image Processing*, Vol. 17, No. 1, 2008.

# 低ランク行列推定による辞書学習

スペースコーディング

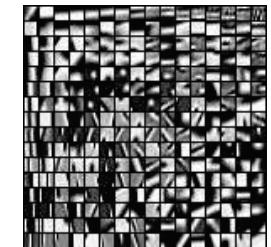
$$\text{観測画像} = \text{イメージパッチ} + \text{ノイズ}$$

$$x = D\alpha + \xi$$

辞書  
(イメージパッチ)  
観測画像  
スペースな係数  
ノイズ  
学習された辞書

$$4.23 \quad 0 \quad 1.24 \quad 0 \quad 0$$

$$\text{イメージパッチ}_1 + \text{イメージパッチ}_2 + \text{イメージパッチ}_3 + \text{イメージパッチ}_4 + \text{イメージパッチ}_5 + \dots$$



$$(\hat{D}, \hat{\alpha}) = \arg \min_{D \in \mathbb{R}^{p \times k}, \alpha_i \in \mathbb{R}^k} \frac{1}{n} \sum_{i=1}^n \|x_i - D\alpha_i\|^2 + \lambda \sum_{i=1}^n \|\alpha_i\|_1$$

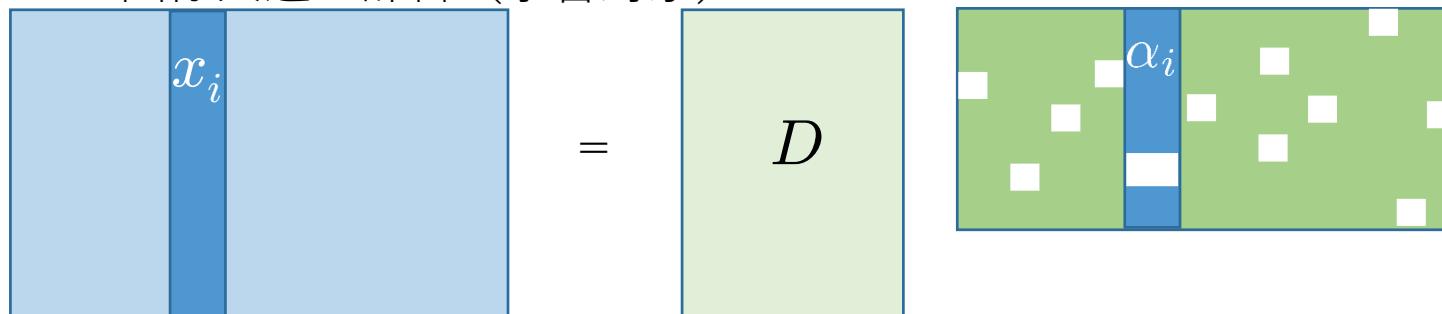
s.t.  $\|D_{:,j}\| \leq 1 \quad (j = 1, \dots, k)$

各画像が  
スペースな係数  
で表現できるよ  
うに  
辞書を構成

$x_i \quad (i=1, \dots, n)$  :  $n$ 枚の画像

$\alpha_i \quad (i=1, \dots, n)$  :  $n$ 枚の画像それぞれの係数 (学習対象)

$D$  : 全画像共通の辞書 (学習対象)



実際はイメージパッチ ( $D$ ) と係数 ( $\alpha$ ) を交互に最適化して学習.

# スペースコーディングを用いたデノイジング

97



This image is taken from MLSS2012 tutorial by F. Bach.

Mairal et al.: Non-local sparse models for image restoration.

In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2009.

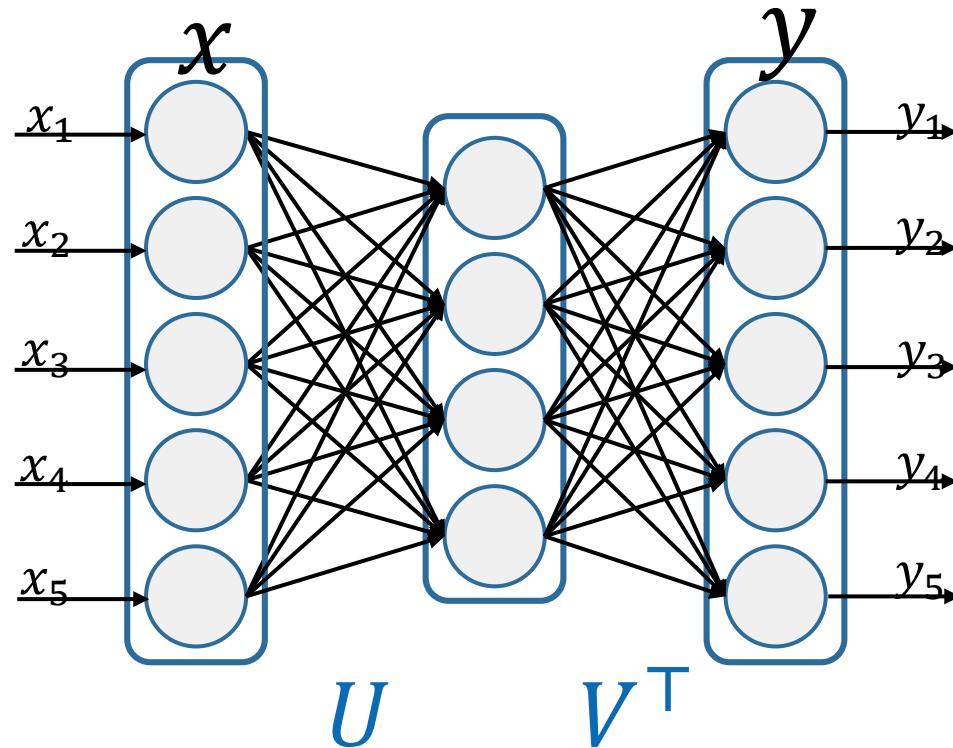
# スペース表現を用いた画像補完

98



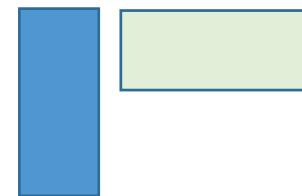
Mairal, Elad and Sapiro: Sparse Representation for Color Image Restoration.  
*IEEE Transactions on Image Processing*, Vol. 17, No. 1, 2008.

# 3層ニューラルネットワーク



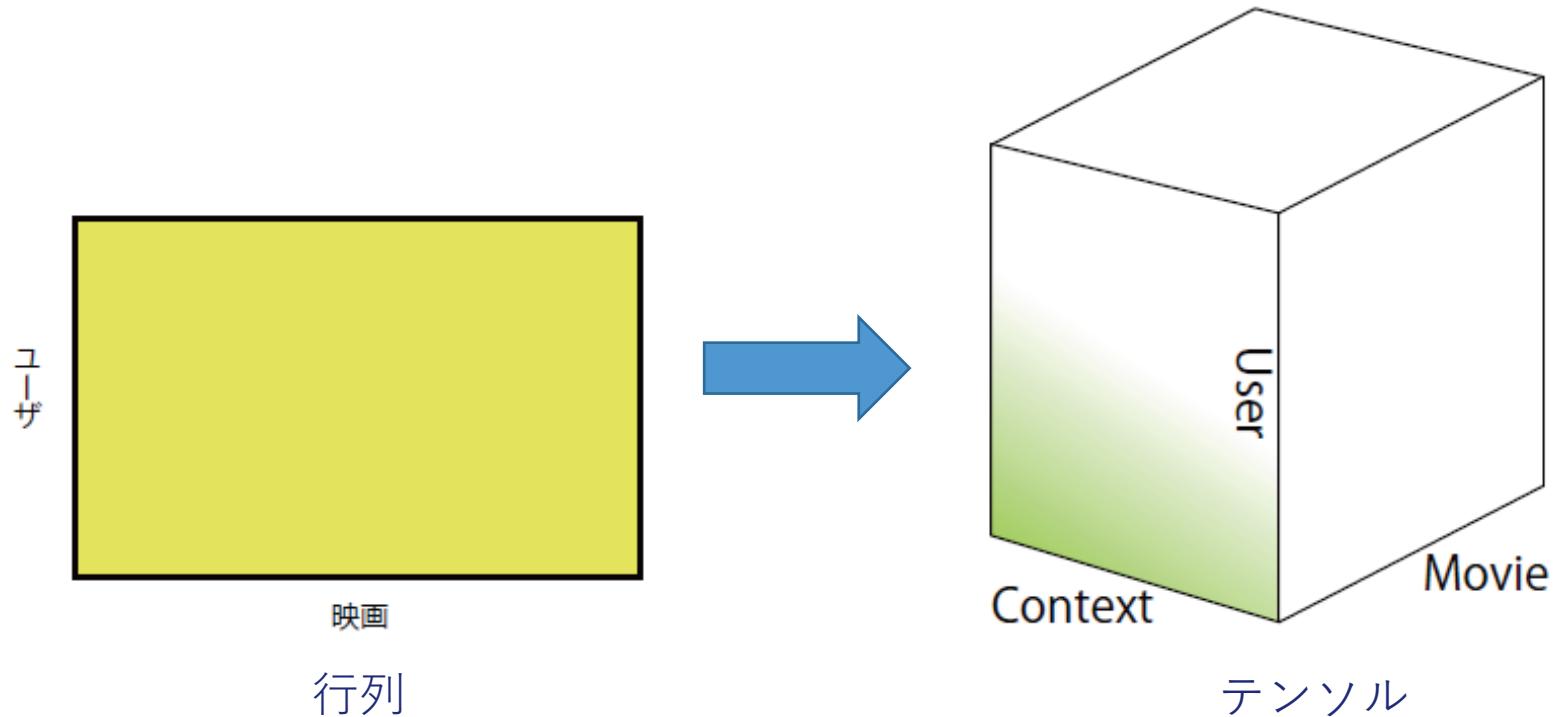
$$f(x) = \underline{V}^T U x \quad \text{低ランク行列}$$

$$f(x) = V^T h(Ux)$$



- 縮小ランク回帰
- マルチタスク学習

# テンソルへの拡張



$$X_{ij} = \sum_{r=1}^d u_{r,i}^{(1)} u_{r,j}^{(2)}$$

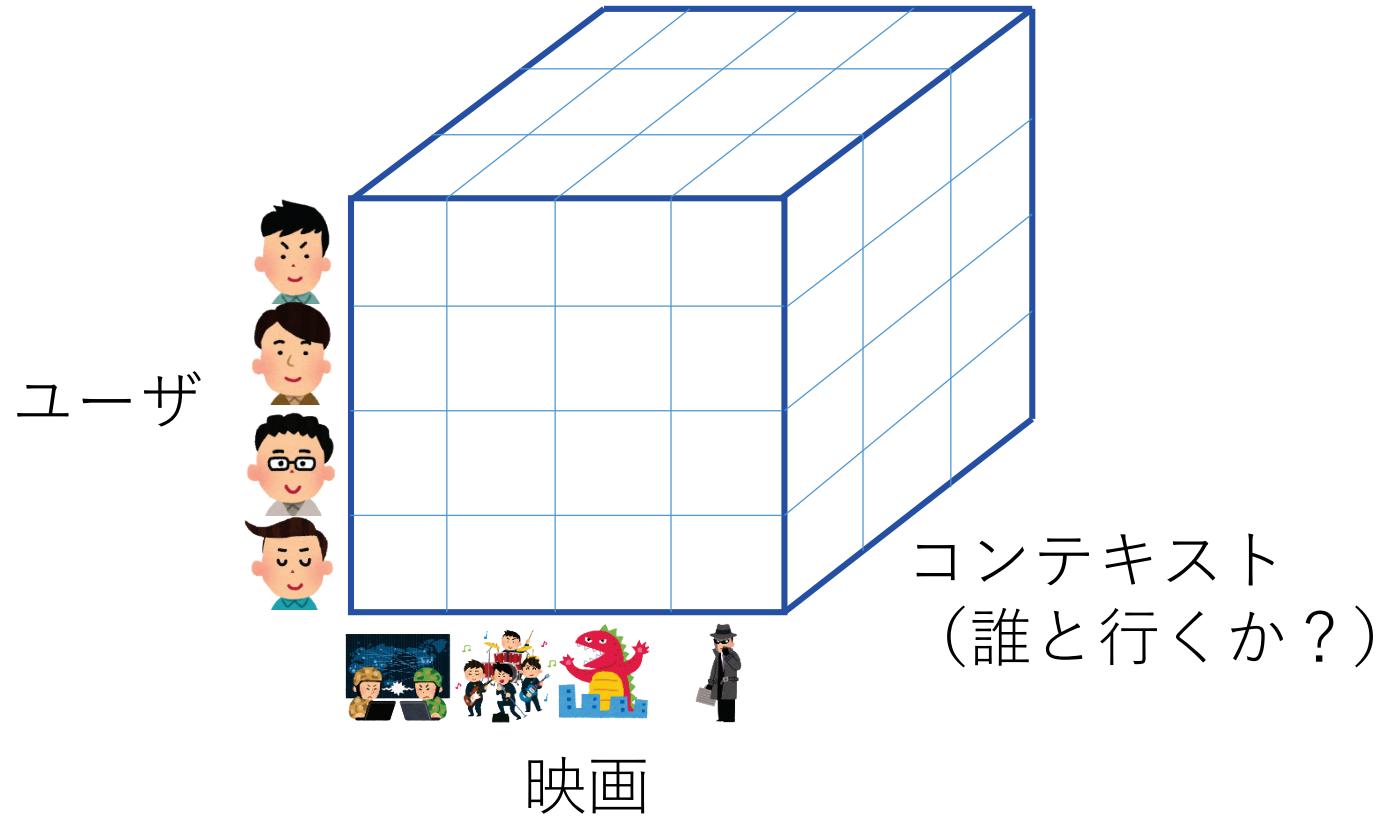
$$X_{ijk} = \sum_{r=1}^d u_{r,i}^{(1)} u_{r,j}^{(2)} u_{r,k}^{(3)}$$

## 応用

- 推薦システム
- 自然言語処理（単語のベクトル表現）
- 時空間データ解析
- 関係データ解析
- マルチタスク学習

# 補助情報も入れた推薦

101

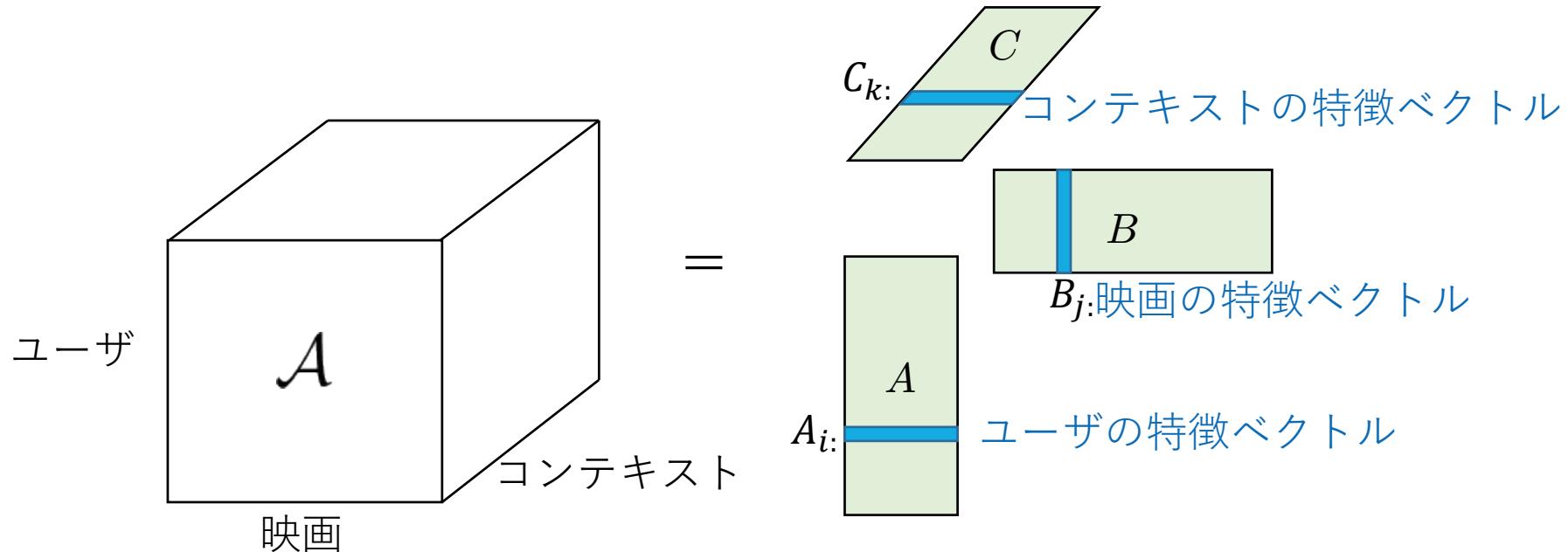


推薦システム：ユーザー × 商品 × 季節

広告モデル：ユーザー × サイト × 広告

バイクシェアリング：会員種類 × 時間 × 位置

# 低ランクテンソルモデル



$$\mathcal{A}_{ijk} = \underbrace{A_{i1}B_{j1}C_{k1}}_{\text{User } i \text{が持つ因子1の重み}} + \underbrace{A_{i2}B_{j2}C_{k2}}_{\text{映画 } j \text{が持つ因子1の重み}} + \cdots + \underbrace{A_{id}B_{jd}C_{kd}}_{\text{コンテキスト } k \text{が持つ因子1の重み}}$$

ユーザ*i*が持つ  
因子1の重み

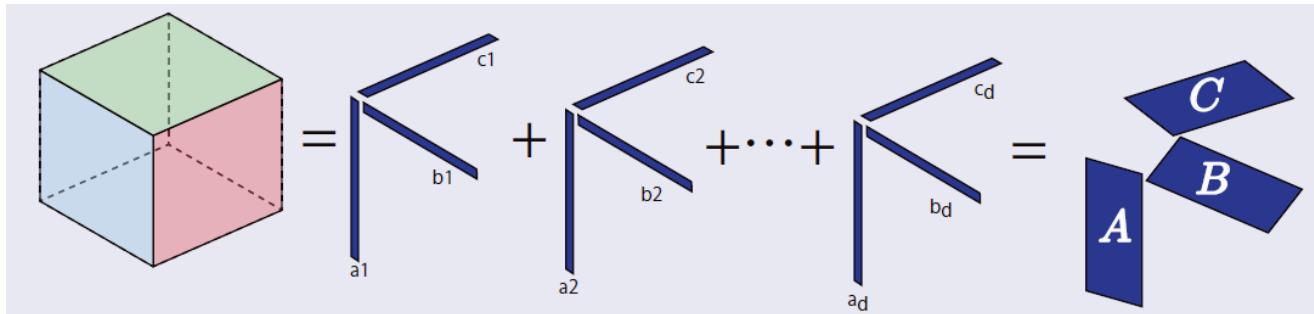
映画*j*が持つ  
因子1の重み

コンテキスト*k*が持つ  
因子1の重み

A,B,Cを観察することで学習結果の解釈も可能

# テンソル分解

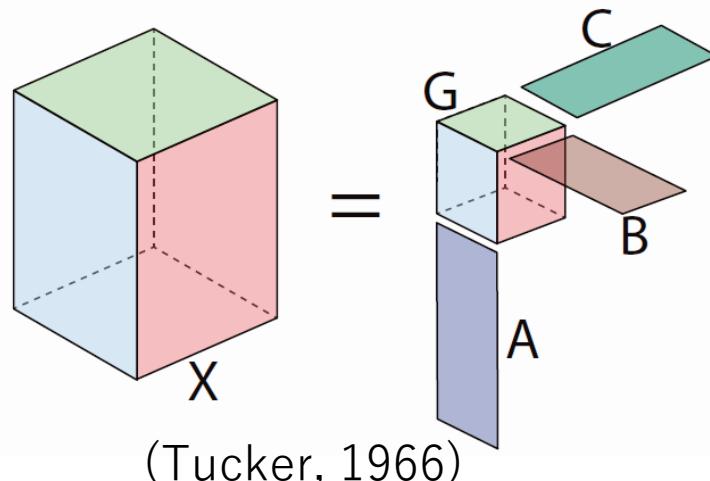
## CP-分解/ランク



Canonical Polyadic 分解(Hitchcock, 1927; Hitchcock, 1927)  
CANDECOMP/PARAFAC (Carroll & Chang, 1970; Harshman, 1970)

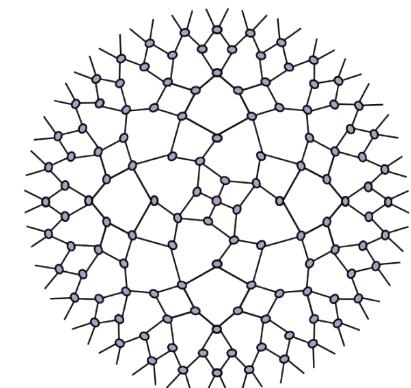
計算はNP困難, ある条件の元で分解の一意性あり

## Tucker-分解/ランク



特異値分解で計算可能

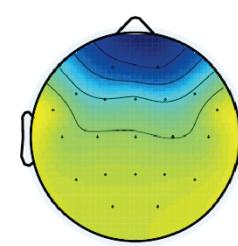
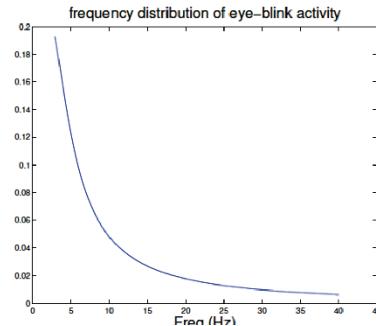
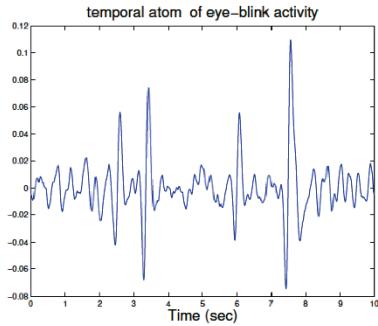
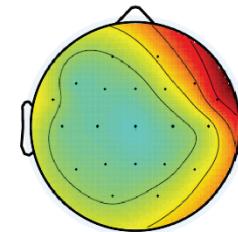
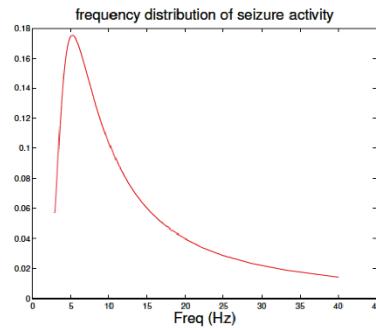
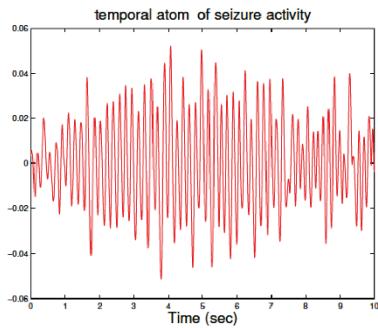
## テンソルネットワーク



(物理学で発展)

# 時空間解析への応用例

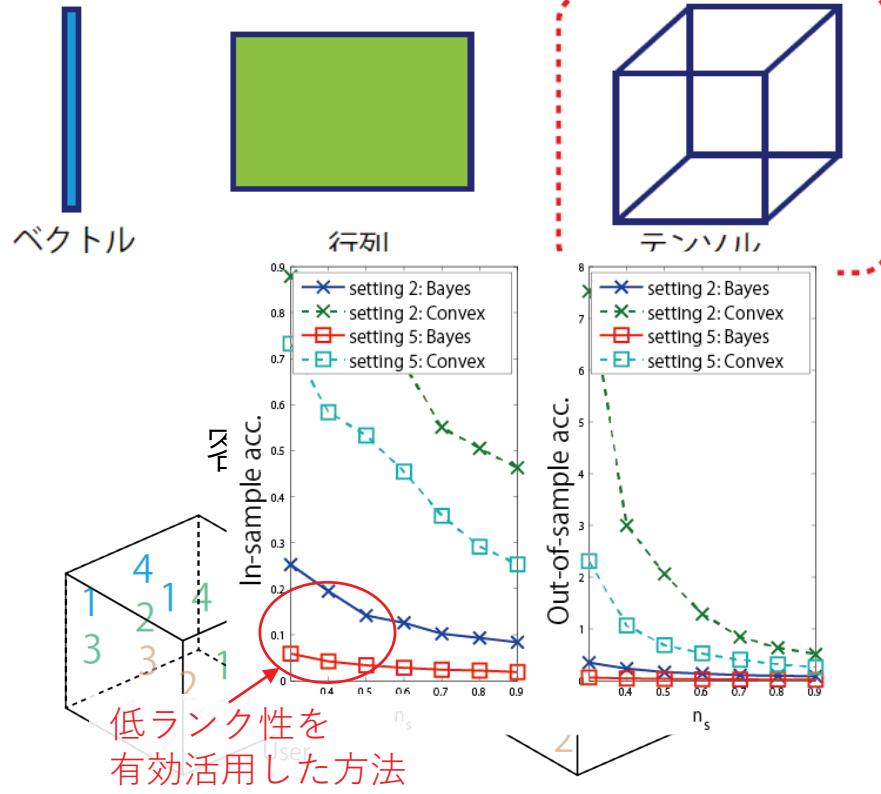
## EEGデータ解析



time × frequency × space  
CP分解

EEG monitoring: Epileptic seizure onset localization (De Vos et al., 2007)

# テンソルの学習



通常(最小二乗法)

$$M^K/n \rightarrow dKM/n$$

次元の呪いを解消

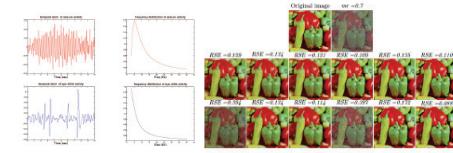


$K$ : 次元  
 $d$ : ランク ( $\ll M$ )

[Suzuki, ICML2015; Kanagawa+et al., ICML2016; Suzuki+et al., NIPS2016]

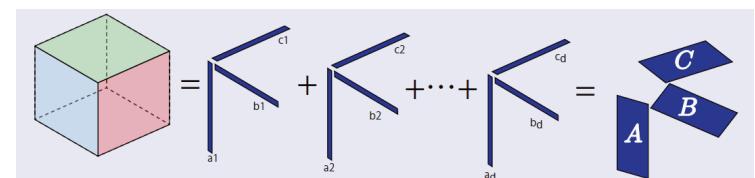
他の応用例：

- 時空間データ解析
- 画像処理
- 自然言語処理
- 深層学習



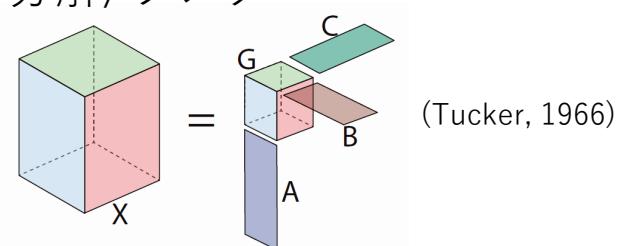
テンソルの“ランク”

CP-分解/ランク



Canonical Polyadic 分解(Hitchcock, 1927; Hitchcock, 1927)  
CANDECOMP/PARAFAC (Carroll & Chang, 1970; Harshman, 1970)

Tucker-分解/ランク



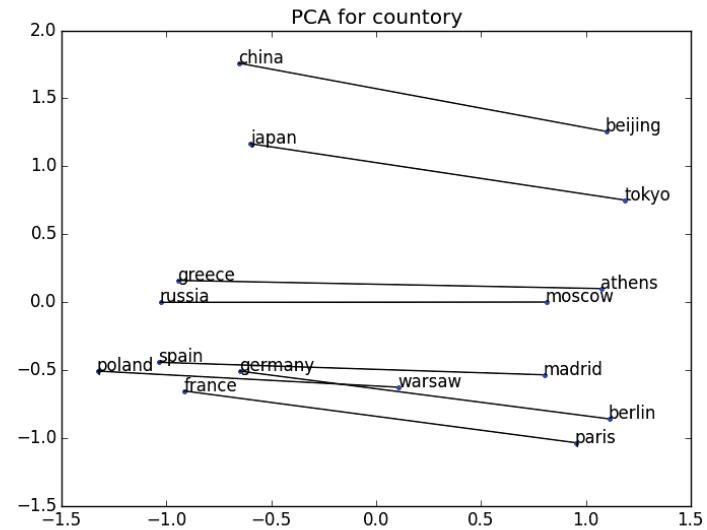
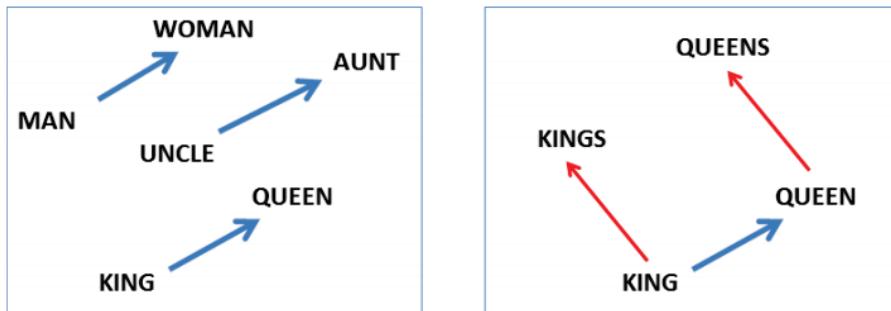
自然言語処理に現れる低ランク性

# Word2vec [Mikolov et al., 2013]

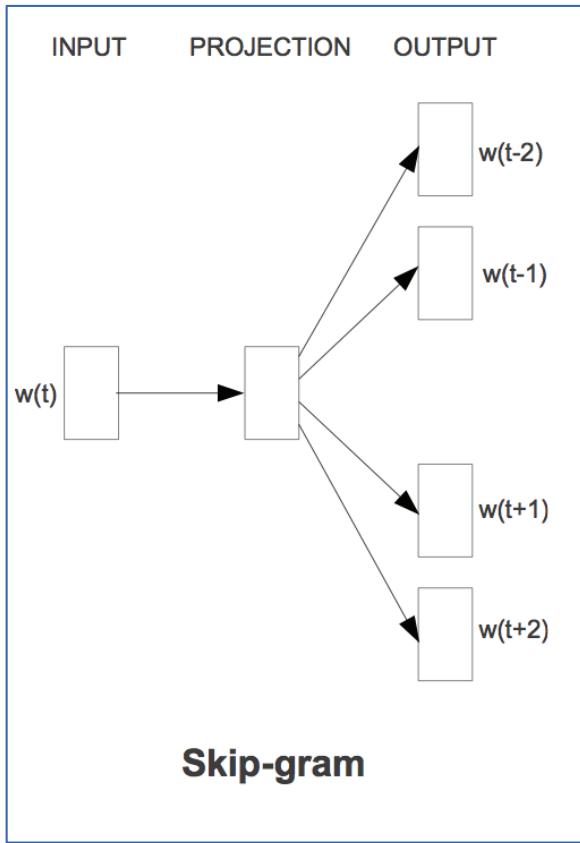
- 単語のベクトル表現を得る方法

“King” – “Man” + “Woman” = “Queen”  
 “Tokyo” – “Japan” + “China” = “Beijing”

意味を足し引きできるような表現が得られる。



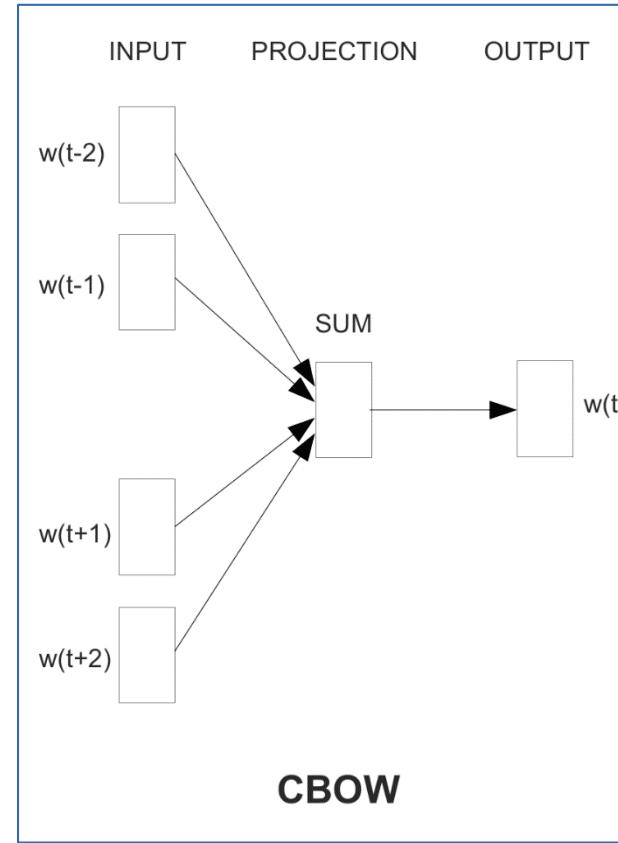
# skip-gramとContinuous Bag-of-Words (CBOW)



Skip-gram

skip-gramモデル

ある単語のまわりに出現する  
単語の確率分布をモデル化



CBOW

CBOWモデル

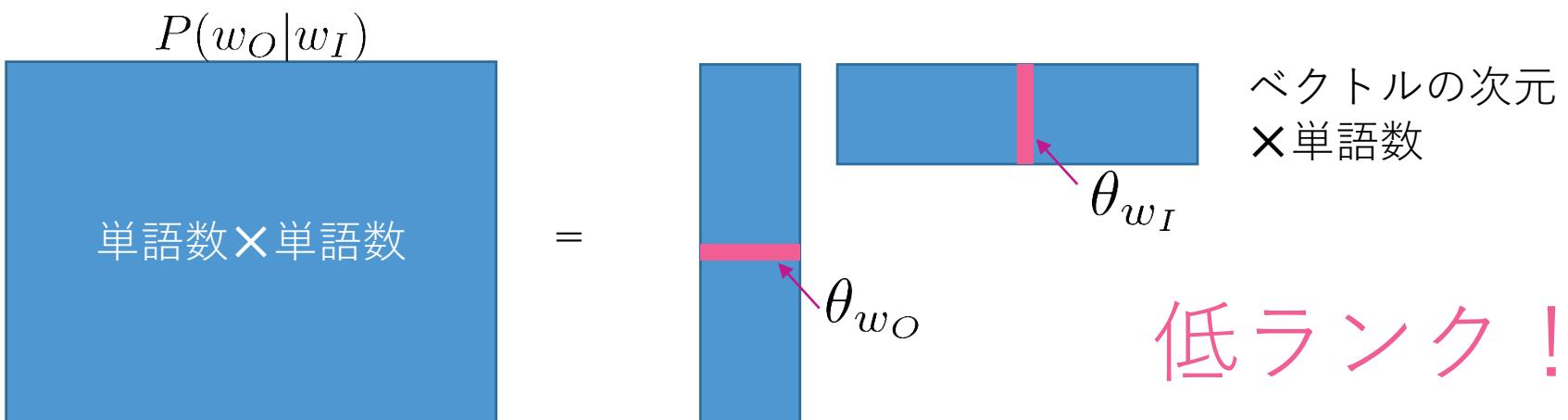
まわりの単語からその場所に  
ある単語が出現する確率をモデル化

# Skip-gramモデル

単語 $w_O$ が $w_I$ の周辺（前後 10 単語ほど）に現れる確率

$$P(w_O|w_I) = \frac{\exp(\langle \theta_{w_O}, \theta_{w_I} \rangle)}{\sum_{w'} \exp(\langle \theta_{w'}, \theta_{w_I} \rangle)}$$
$$\propto \exp(\langle \theta_{w_O}, \theta_{w_I} \rangle)$$

- 単語のベクトル表現  $\theta_{w_O}$  と  $\theta_{w_I}$  の内積で表現.
- ベクトルの次元はせいぜい500ほど



# 実際の挙動

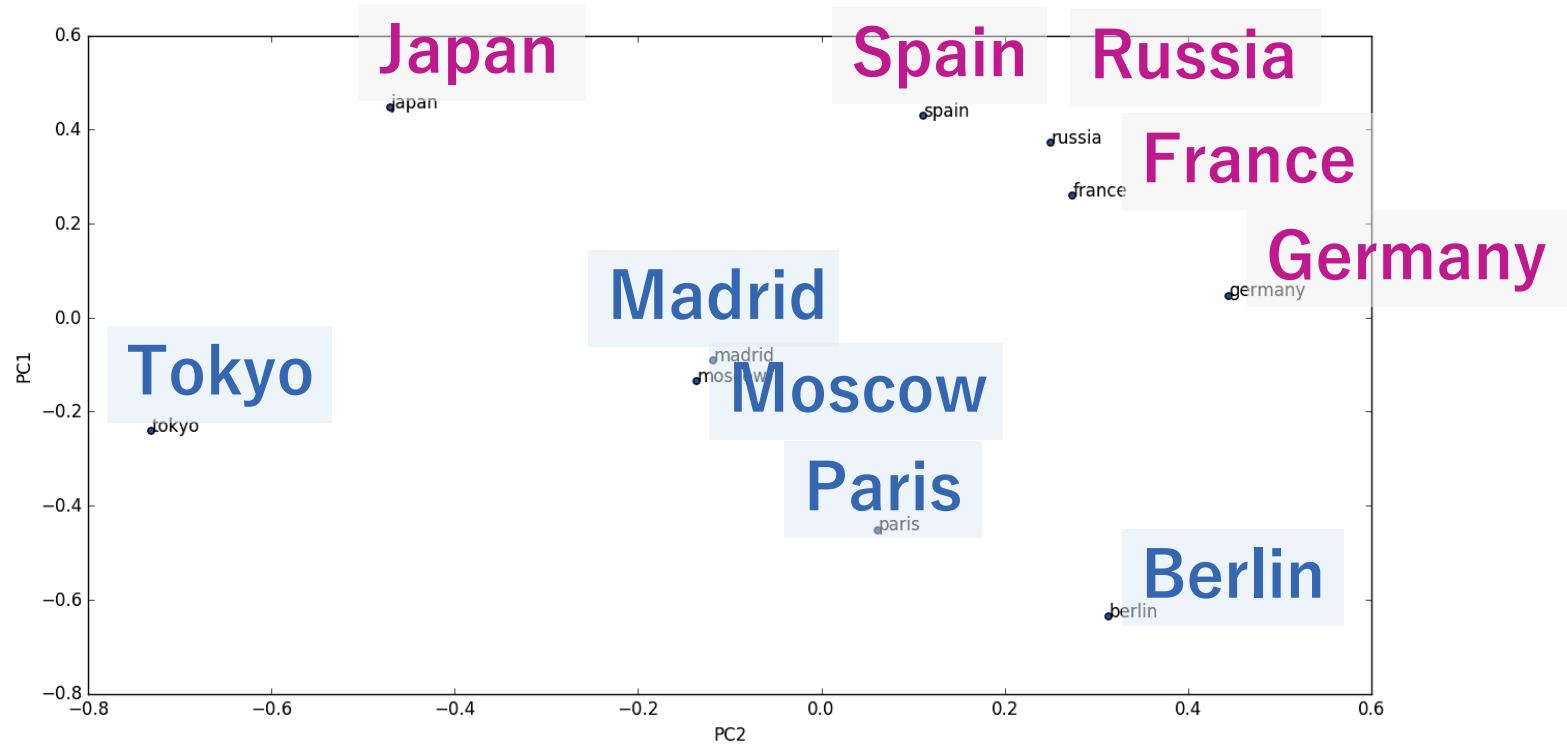
```
from gensim.models import word2vec  
  
train_file = "./mldata/text8"  
  
data = word2vec.Text8Corpus(train_file)  
model = word2vec.Word2Vec(size=100, window=5, min_count=5, workers=7)  
  
model.build_vocab(data)  
model.train(data)
```

次元 100, 前後 5 単語の出現頻度をモデル化, 5 回以下の出現単語は無視

“Queen” + “Man” - “Woman” = “King” ?

```
>>> model.most_similar(positive=['queen','man'],negative=['woman'])  
[('king', 0.6050819158554077), ('scotland', 0.587989091873169), ('prince', 0.573  
6681222915649), ('elizabeth', 0.571208119392395), ('lord', 0.5638244152069092),  
('duchess', 0.5520190000534058), ('duke', 0.5498123168945312), ('crown', 0.54618  
62087249756), ('sir', 0.5441839694976807), ('lorraine', 0.5441141128540039)]
```

# 国と首都



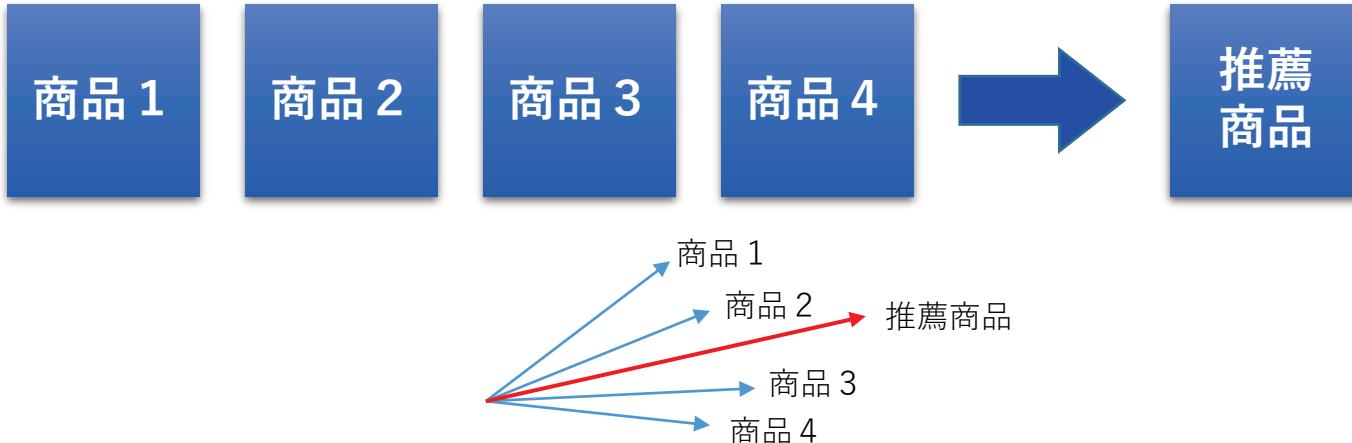
word2vecの貢献：

データの「意味」を低次元ベクトルとして表現できることを実験的に示した。  
→ 深層学習にもつながる考え方。

# Word2vecの応用

- 商品タグから関連した商品を推薦

あるユーザーの購買履歴



- 口コミの要約

## 文章をベクトル化する手法：skip-thought, text2vec

- 記事から関連広告を推薦
- 感情分析

# まとめ

- 機械学習の歴史
- 機械学習の考え方
  - 複雑な規則をデータから学ぶ
- モデルと損失
  - 学習：期待損失最小化
- 過学習の問題
  - 複雑なモデルを当てはめれば良いわけではない.
  - 正則化
  - 変数選択
- スペース高次元データ解析
  - Lasso
- 低ランク行列/テンソル分解  
→深層学習にもつながる

# 補足資料

# 最適化法

$$\min_{\beta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \ell(y_i, x_i^\top \beta) + \psi(\beta)$$

手元にあるデータへの当てはまり  
(損失関数)

正則化項

記法を簡略化

$$\min_{\beta \in \mathbb{R}^p} f(\beta) + \psi(\beta)$$

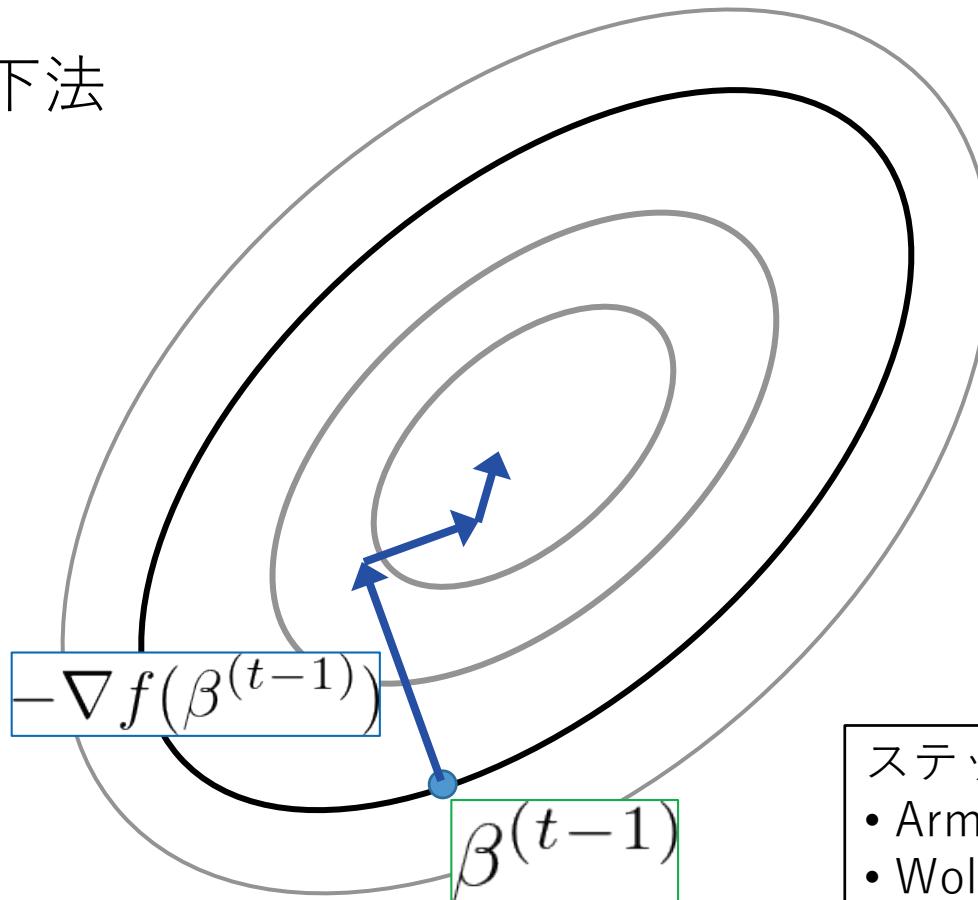


最適化法：どうやってこの最適化問題を解く？

- 勾配法
- 座標降下法
- 交互方向乗数法

# 正則化項がない場合の最適化

- 最急降下法



ステップサイズの決定には  
• Armijoの規準  
• Wolfeの規準  
等がある。

$$\beta^{(t)} = \beta^{(t-1)} - \alpha_t \nabla f(\beta^{(t-1)})$$

# 正則化項がある場合：近接勾配法

$$f(\beta) + \psi(\beta)$$

↓ 線形近似      ↓ 正則化項はそのまま  
(正則化項は微分不可能)

$$g_t = \nabla f(\beta^{(t-1)})$$

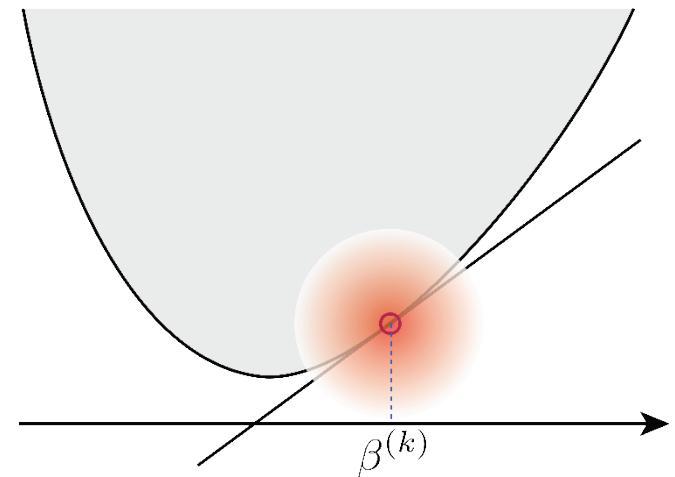
$$\beta^{(t)} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ g_t^\top \beta + \psi(\beta) + \frac{\eta_t}{2} \|\beta - \beta^{(t-1)}\|^2 \right\}$$

↑ 遠くへ離れないようにする項

鍵となる計算：近接写像

$$\text{prox}(\mathbf{q}|\psi) := \arg \min_{\mathbf{x}} \left\{ \psi(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{q}\|^2 \right\}$$

→ L1正則化なら簡単に計算可能。  
(Soft-thresholding関数)



# 近接写像を用いた表記

$$g_t = \nabla f(\beta^{(t-1)})$$

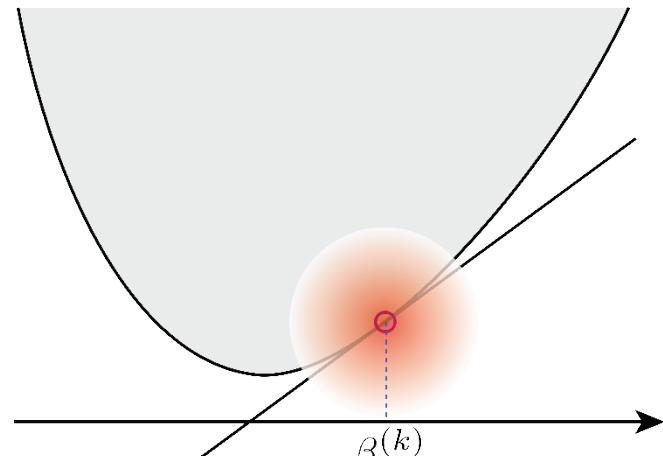
$$\begin{aligned}\beta^{(t)} &= \arg \min_{\beta \in \mathbb{R}^p} \left\{ g_t^\top \beta + \psi(\beta) + \frac{\eta_t}{2} \|\beta - \beta^{(t-1)}\|^2 \right\} \\ &= \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{\psi(\beta)}{\eta_t} + \frac{1}{2} \left\| \beta - \left( \beta^{(t-1)} - \frac{g_t}{\eta_t} \right) \right\|^2 \right\} \\ &= \text{prox}(\beta^{(t-1)} - g_t / \eta_t | \psi / \eta_t)\end{aligned}$$

- $\Psi=0$ なら単なる最急降下法
- $\Psi$ がある凸集合の標示関数なら射影勾配法

鍵となる計算：近接写像

$$\text{prox}(\mathbf{q} | \psi) := \arg \min_{\mathbf{x}} \left\{ \psi(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{q}\|^2 \right\}$$

→ L1正則化なら簡単に計算可能。  
(Soft-thresholding関数)



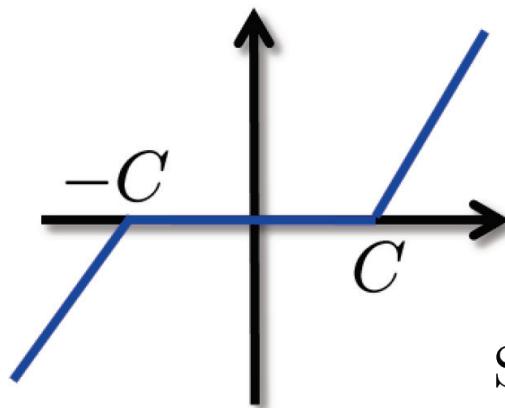
# 具体例：L1正則化

座標ごとにわかっている！

$$\beta_j^{(t)} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ g_{t,j} + \lambda |\beta_j| + \frac{\eta_t}{2} (\beta_j - \beta_j^{(t-1)})^2 \right\}$$

$$\begin{aligned} \text{prox}(\mathbf{q} | C \| \cdot \|_1) &= \arg \min_{\mathbf{x}} \left\{ C \|\mathbf{x}\|_1 + \frac{1}{2} \|\mathbf{x} - \mathbf{q}\|^2 \right\} \\ &= (\underbrace{\text{sign}(q_j) \max(|q_j| - C, 0)}_{\text{ST}_C(q_j) \text{ とおく}})_j. \end{aligned}$$

$$\beta_j^{(t)} = \text{ST}_{\lambda/\eta_t} \left( \beta_j^{(t-1)} - \frac{g_{t,j}}{\eta_t} \right) \quad \text{解が陽に書ける！}$$



# 近接勾配法の収束速度

$$\beta^{(t)} = \text{prox}(\beta^{(t-1)} - g_t/\eta_t | \psi/\eta_t)$$

$f$ の性質	$\mu$ -強凸	非強凸
$L$ -平滑	$\exp\left(-t\frac{\mu}{L}\right)$	$\frac{L}{t}$
非平滑	$\frac{1}{\mu t}$	$\frac{1}{\sqrt{t}}$

- 滑らかなほど・強凸なほど速い。
- Nesterov の加速法を用いれば滑らかな場合に速くなる (Nesterov, 2007, Zhang et al., 2010).
- 上のオーダーは勾配情報のみを用いる方法 (First order method) の中で最適。

$\eta_t$ の設定	強凸	非強凸
平滑	$L$	$L$
非平滑	$\frac{\mu t}{2}$	$\sqrt{t}$

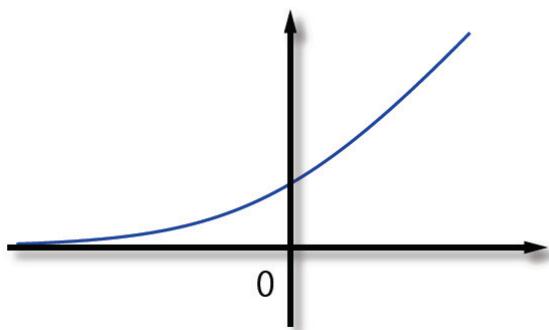
# 強凸性と平滑性

- 平滑性：勾配の変化がゆっくり

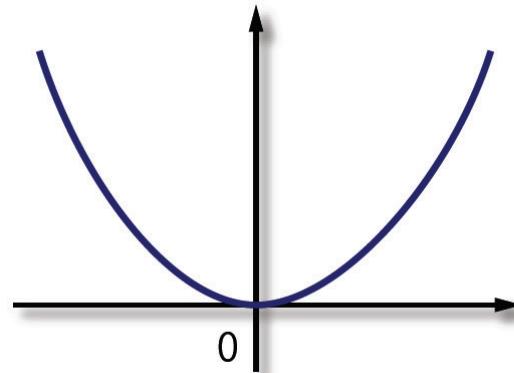
$$\|\nabla f(x) - \nabla f(x')\| \leq L \|x - x'\|$$

- 強凸性：2次関数以上に曲がっている

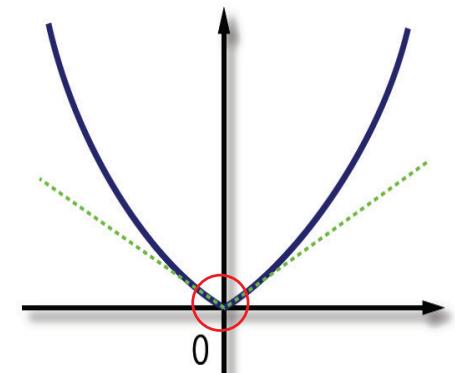
$$\frac{\mu}{2} \theta(1-\theta) \|x-y\|^2 + f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y)$$



平滑だが強凸ではない



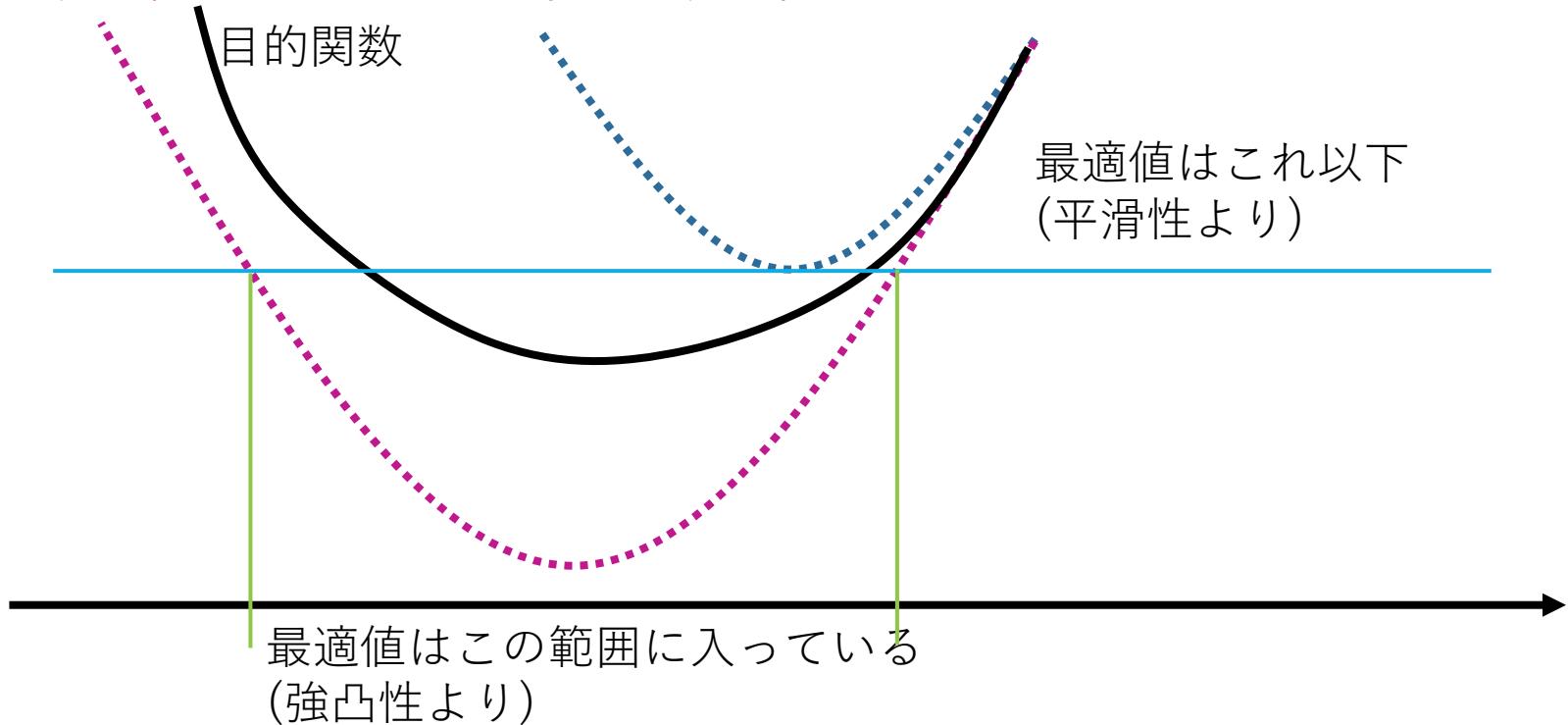
平滑かつ強凸



強凸だが平滑ではない

強凸性による  
下からのバウンド

平滑性による  
上からのバウンド



平滑性 → 最適値を上から抑えられる.  
強凸性 → 最適値の範囲を限定できる.

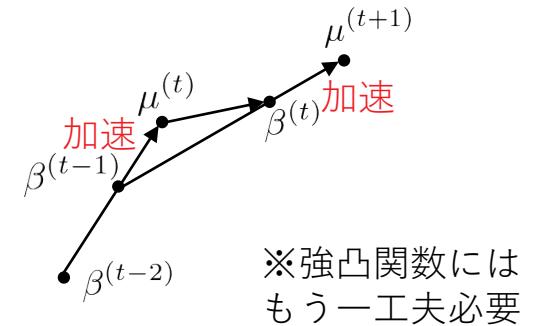
# 近接勾配法の収束速度

## Nesterovの加速法

$$\beta^{(t)} = \text{prox}(\mu^{(t)} - g_t / \eta | \psi / \eta)$$

$$s_{t+1} = \frac{1 + \sqrt{1 + 4s_t^2}}{2}$$

$$\mu^{(t+1)} = \beta^{(t)} + \left( \frac{s_t - 1}{s_{t+1}} \right) (\beta^{(t)} - \beta^{(t-1)})$$

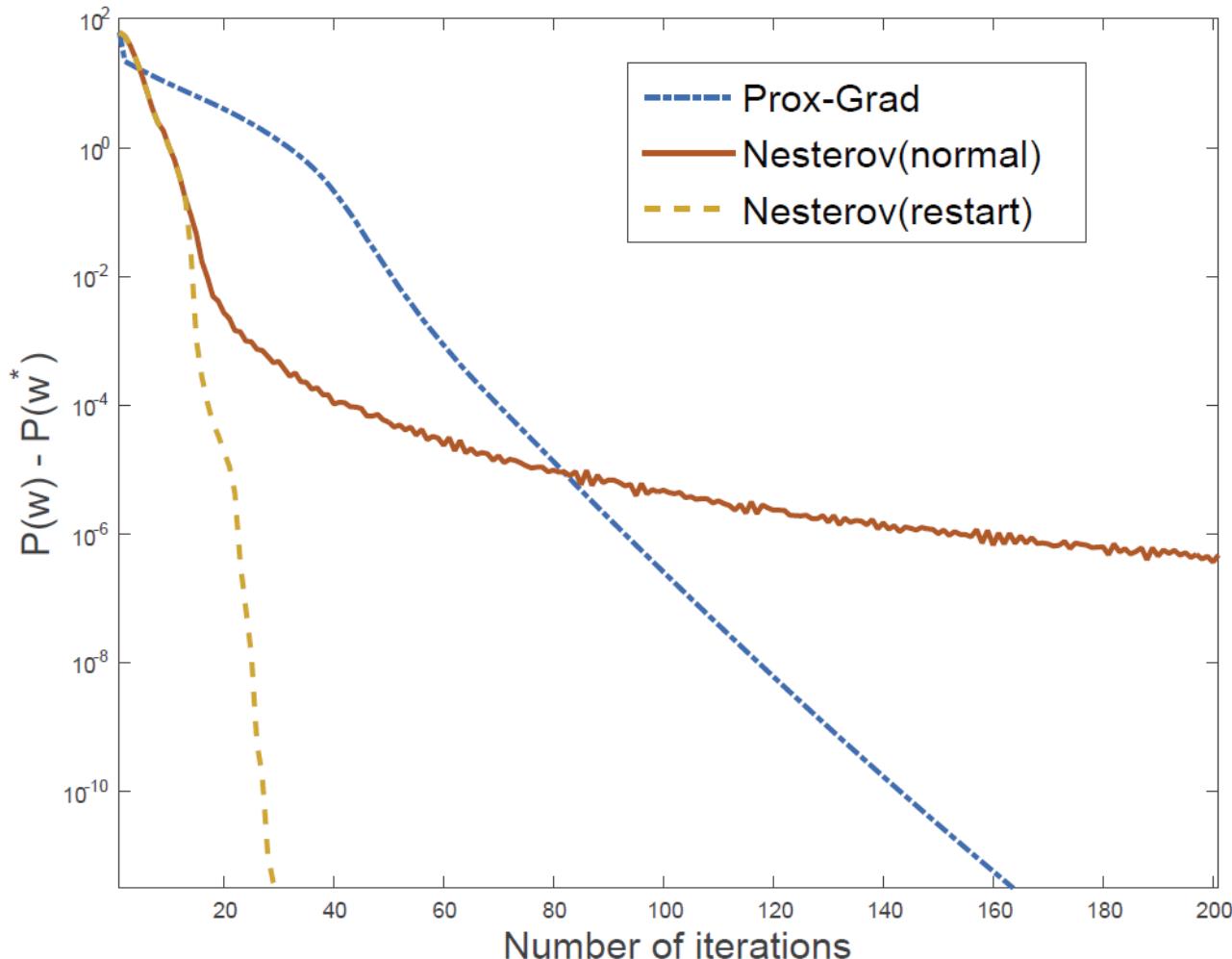


$f$ の性質	$\mu$ -強凸	非強凸
$L$ -平滑	$\exp\left(-t\sqrt{\frac{\mu}{L}}\right)$	$\frac{L}{t^2}$
非平滑	$\frac{1}{\mu t}$	$\frac{1}{\sqrt{t}}$

- 滑らかなほど・強凸なほど速い。
- Nesterov の加速法を用いれば滑らかな場合に速くなる (Nesterov, 2007, Zhang et al., 2010).
- 上のオーダーは勾配情報のみを用いる方法 (First order method) の中で最適。

$\eta_t$ の設定	強凸	非強凸
平滑	$L$	$L$
非平滑	$\frac{\mu t}{2}$	$\sqrt{t}$

# 数値実験例



サンプルサイズ  $n = 700$ , 次元  $p = 1000$ , 100成分のみ非ゼロ

# 総計算量

経験損失は  $O(1/n)$  まで下げる必要がある。 (汎化誤差ミニマックスレートが  $O(1/n)$ )

$$\kappa = \frac{L}{\mu} : \text{条件数}$$

一回の更新にかかる計算量

$$O(n)$$

$\times$

$1/n$  まで下げるのに必要な更新回数

$$O(\sqrt{\kappa} \log(n))$$

$$\nabla f(\beta) = \sum_{i=1}^n \nabla_{\beta} \ell(z_i, \beta) : n \text{個の和}$$

Nesterovの加速法

=

総計算量

$$O(n \sqrt{\kappa} \log(n))$$

- サンプルサイズ  $n$  が強く効いてくる
  - 大規模データの最適化が難しい
- 確率的最適化が有用

