

機械学習のお話

LTの目標

- 機械学習の全体像をざっくり把握する. (統計知らない人)
- 機械学習を統計学的視点からざっくり掴む.
(統計知ってる人)

機械学習とは

経験 E , タスク T , パフォーマンス尺度 P

「(機械) 学習とは」 :

タスク T について, P で測られるタスクの実行能力が
経験 E を通じて向上すること.

例) 手書き文字認識(T) において,
書いてある数字を当てる正答率(P)が
データ(E) を通じて向上する.

最近流行りのワード

よく分からない



教師あり学習

分類問題

ベイジアンネットワーク

クラスタリング

SVM

深層学習

学習法での分類

- 教師あり学習
- 教師なし学習
- 強化学習
- 半教師つき学習
- 転移学習
- マルチタスク学習

教師あり学習 (supervised learning)

学習データに正解ラベルを付けて学習する方法.

訓練データ

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

x が訓練用のデータベクトル. y が正解の値.

x から正しい y の値を出力する関数を作る.

(例.)

アヤメのデータからアヤメの種類を分類する.

$x_i = (\text{がく片の長さ}, \text{がく片の幅}, \text{花びらの長さ}, \text{花びらの幅})$

$y_i = (0: \text{ヒオウギアヤメ}, 1: \text{ブルーフラッグ}, 2: \text{Virginica})$

新しい x_* を与えたら 正しく種類が判定される関数を作る.

ラベル	がく片の 長さ	がく片の 幅	花びらの 長さ	花びらの 幅
ヒオウギ アヤメ	3.0	1.0	5.0	10.0
ヒオウギ アヤメ	3.4	1.3	4.4	1.1

⋮

教師なし学習

学習データに正解ラベルを与えず,
データのみからデータの本質的構造を抽出する学習法.

訓練データ

$$D = \{ x_1, x_2, \dots, x_n \}$$

(例) : 生徒のテストデータから良いクラスタリングを見つける.
100人の生徒のテスト結果

$x_i = (\text{国語の点}, \text{数学の点}, \text{英語の点}, \text{社会の点}, \text{理科の点})$

生徒をクラス内で実力差の少ない3つのクラスに分ける.

国語	数学	理科	社会	英語

強化学習

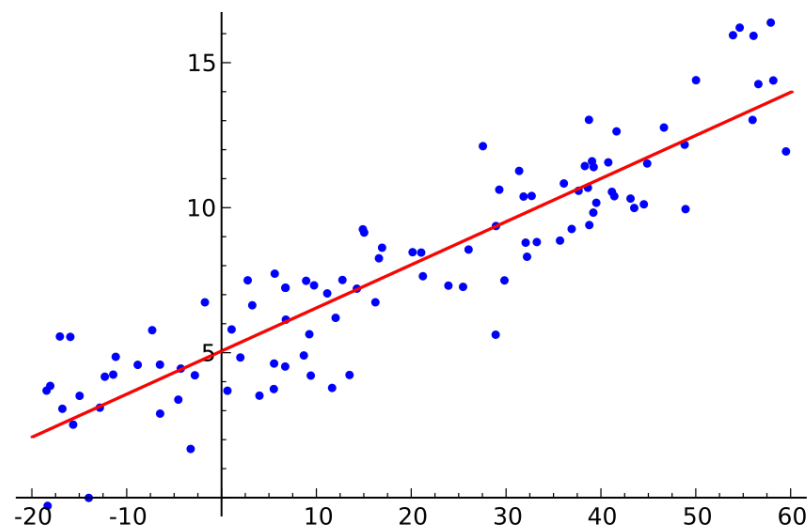
- ごめんよく知りません。

モデルでの分類

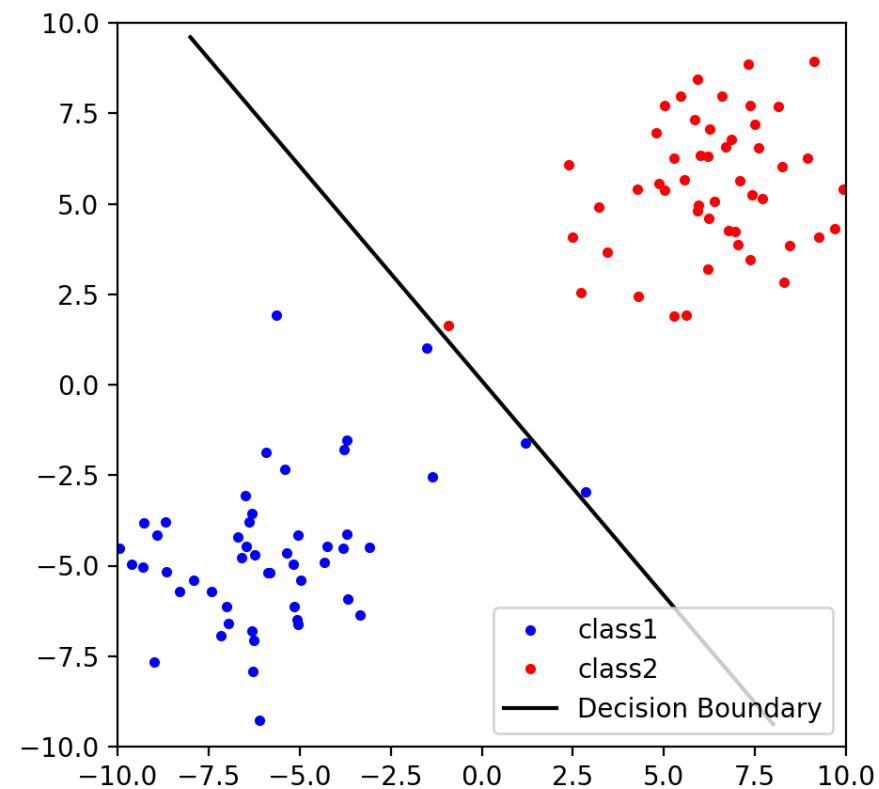
- 線形モデル
- 加法モデル (+正則化)
- カーネルモデル
 (線形モデルと組み合わせるカーネルトリック. SVMは有名)
- 深層モデル

他に木やブースティングなどのモデルもある。

線形モデル



線形回帰



線形分類

深層モデル

<https://qiita.com/hiyoko9t/items/089c3c828b2d3b04f756>

「人工知能でバレンタインチョコが本命か義理かを判別する」

•本命の確率が90%以上と予測されているもの



•義理の確率が90%以上と予測されているもの



CNN (Convolution Neural Network) : 画像認識分野で非常に強力

学習法とモデルの関係

今流行ってるやつ

	線形	加法	カーネル	深層
教師あり	一般化線形モデル SVM	平滑化スプライン	カーネル密度	CNN, RNN
教師なし	PCA K-means			AE, VAE
強化学習				DQN

易



難

タスク別の分類

学習法による分類をさらに細分化する.

教師あり学習：回帰, 分類

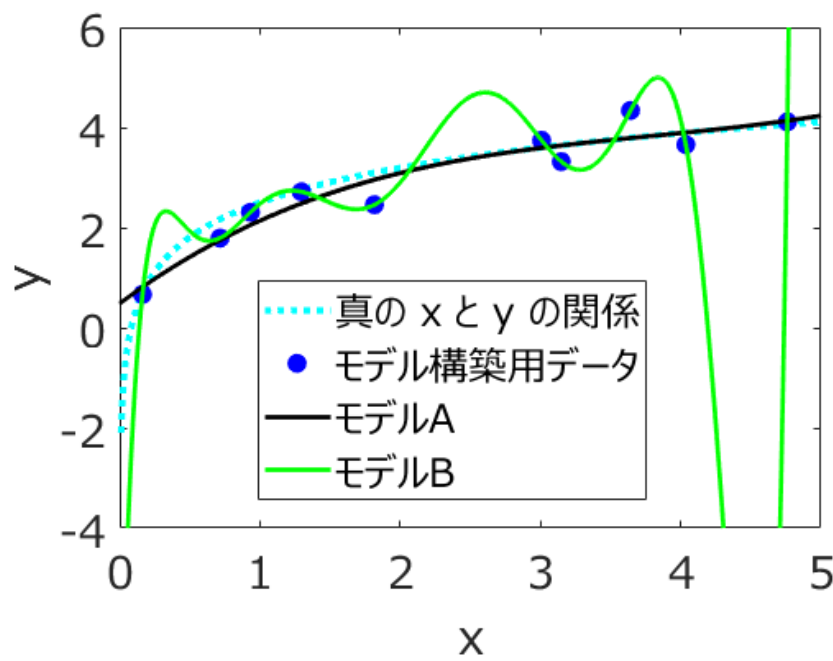
教師なし学習：クラスタリング, 次元圧縮, 密度推定

回帰 (Regression)

データから

$$\boldsymbol{x} \in \mathbb{R}^M, y \in \mathbb{R}, y = f(\boldsymbol{x})$$

を満たす関数 f を発見する. y は連続値をとる.



適切なモデルを作らないと過学習

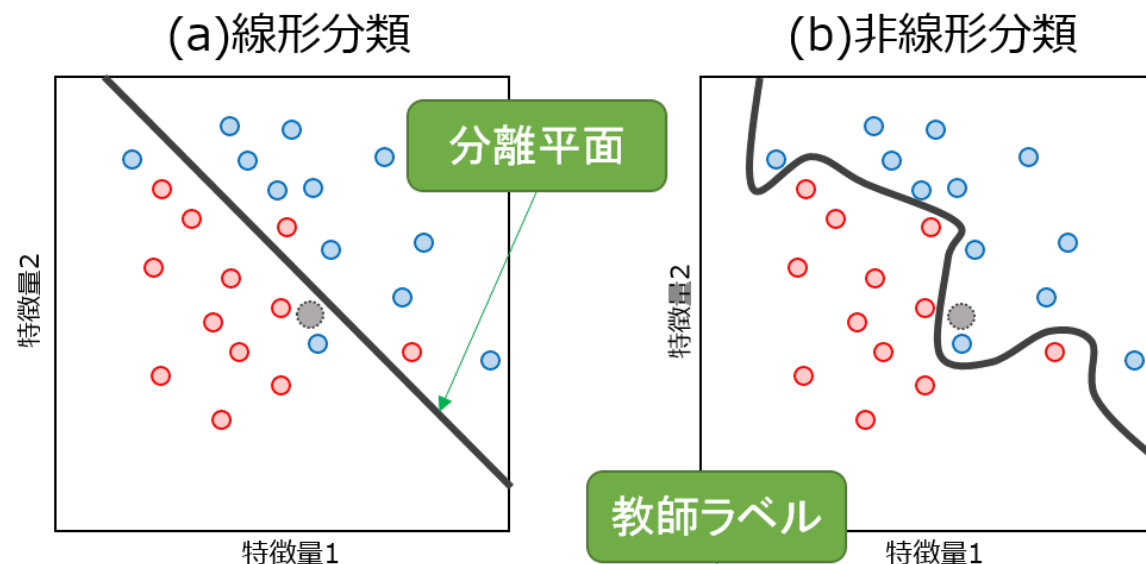
分類 (Classification)

データから

$$\boldsymbol{x} \in \mathbb{R}^M, y \in \{1, 2, \dots, n\}, y = f(\boldsymbol{x})$$

を満たす関数 f を発見する. y は離散値をとる.

例 :

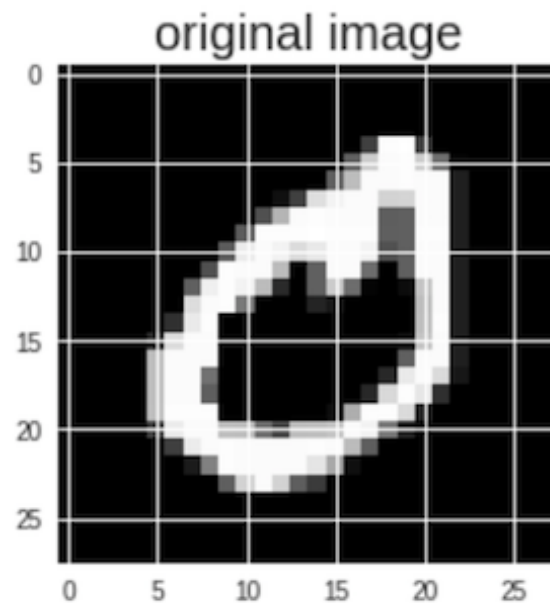


- 学習用データ(教師データ) ● クラス1 ● クラス2
- 分類したいデータ : ● (a)ではクラス1 (b)ではクラス2に分類

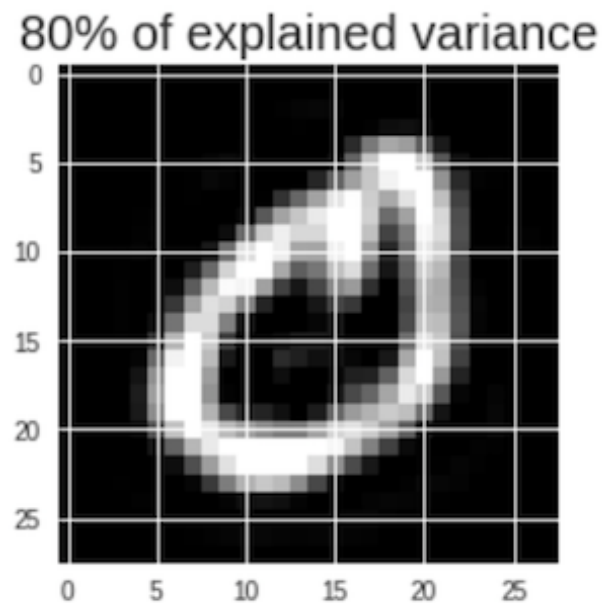
次元削減（ Dimensionality reduction ）

データの次元数を減らし、データの圧縮や潜在構造の分析を行う。

例： PCA（主成分分析）



784個の変数



150 個の変数

統計的機械学習

「機械学習のうち、データの確率的な生成規則を学習するもの」

統計と機械学習の違いって？（諸説あり）

1. 統計学はデータを説明することを重視
2. 機械学習はデータから予測を行うことを重視

確率を用いた機械学習
(統計を知ってる人向け)

確率を用いる機械学習

x : データ

y : (ex.) 回帰曲線の値, 属する分類, 属するクラス

$$\mathbf{x} \in \mathbb{R}^M, y = f(\mathbf{x})$$

- 関数モデル (確率を考えない)
- 機械学習に対する確率的アプローチ

1. 識別モデル

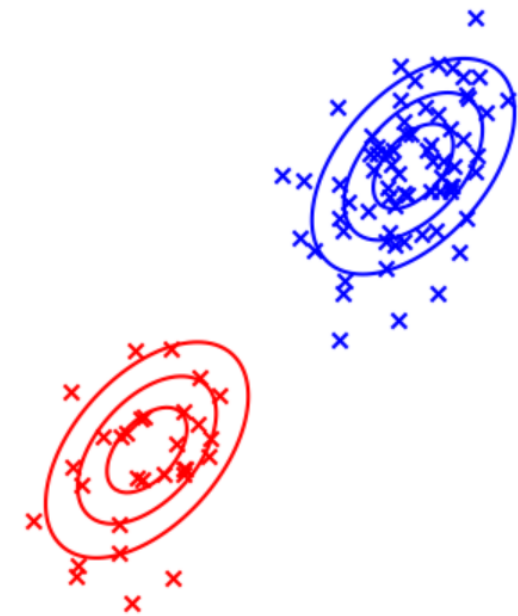
事後確率 $P(y | \mathbf{x})$ を直接モデル化

2. 生成モデル

$P(\mathbf{x} | y)$ と $P(y)$ をモデル化.

$$P(y | \mathbf{x}) = \frac{P(\mathbf{x} | y)P(y)}{P(\mathbf{x})}$$

ベイズの公式を使う.



単回帰（関数モデル）

問題 データ $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ から

$$y = ax + b$$

$a, b \in \mathbb{R}$ はパラメータ. 良い a, b を見つける.

どういう関数を定義したら良さそうか？

$$Error(\mathbf{x}) = \sum_{i=1}^N \{y_i - (ax_i + b)\}^2$$

$$Error(\mathbf{x}) = \sum_{i=1}^N \{y_i - (ax_i + b)\}^4 \quad \text{とか採用しても一見良さそう}$$

単回帰の確率的考察

(準備1) ベクトル表現に直す.

$$\mathbf{x}_i = \begin{bmatrix} 1 \\ x_i \end{bmatrix}, \mathbf{w} = \begin{bmatrix} b \\ a \end{bmatrix}$$

$$\text{minimize} \sum_{i=1}^N \{y_i - (ax_i + b)\}^2 \Leftrightarrow \text{minimize} \sum_{i=1}^N y_i - \mathbf{w}^T \mathbf{x}_i$$

単回帰の確率的考察

(準備2) 最小二乗法を最尤推定として解釈する.

$$\text{minimize } \sum_{i=1}^N (y_i - \boldsymbol{w}^T \boldsymbol{x}_i)^2 \Leftrightarrow \text{maximize } \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi\sigma^2}} \exp \frac{-(y_i - \boldsymbol{w}^T \boldsymbol{x}_i)^2}{2\sigma^2}$$

σ は何かしらの値。

y が正規分布に従っているとした時の

平均の最尤推定量が $\boldsymbol{w}^T \boldsymbol{x}$

単回帰の識別モデルアプローチ

(識別モデルとしての定式化)

仮定:

$$P(y_n \mid \mathbf{x}_n, \mathbf{w}) = N(y_n \mid \mathbf{w}^T \mathbf{x}_n, (\sigma^2)^{-1})$$

$$P(\mathbf{w}) = N(\mathbf{w} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1})$$

$\boldsymbol{\mu}, \boldsymbol{\Sigma}$ も適当に仮定する.

パラメータ $\mathbf{w} = \begin{bmatrix} b \\ a \end{bmatrix}$ の分布を考える発想はベイズ統計学

パラメータ \boldsymbol{w} をデータから推定.

$$\begin{aligned}P(\boldsymbol{w} \mid \boldsymbol{y}, \boldsymbol{x}) &= \frac{P(\boldsymbol{w})P(\boldsymbol{y} \mid \boldsymbol{x}, \boldsymbol{w})}{P(\boldsymbol{y} \mid \boldsymbol{x})} \\&\propto P(\boldsymbol{w})P(\boldsymbol{y} \mid \boldsymbol{x}, \boldsymbol{w}) \\&= N(\boldsymbol{w} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}) \prod_{i=1}^N N(y_i \mid \boldsymbol{w}^T \boldsymbol{x}_i, (\sigma^2)^{-1}) \\&= \dots = N(\boldsymbol{w} \mid \hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}^{-1})\end{aligned}$$

$$\text{ただし, } \hat{\boldsymbol{\Sigma}} = (\sigma)^2 \sum_{i=1}^N \boldsymbol{x}_i \boldsymbol{x}_i^T + \boldsymbol{\Sigma}, \quad \hat{\boldsymbol{\mu}} = \hat{\boldsymbol{\Sigma}}^{-1} \{ (\sigma)^2 \sum_{i=1}^N y_i \boldsymbol{x}_i + \boldsymbol{\Sigma} \boldsymbol{\mu} \}$$

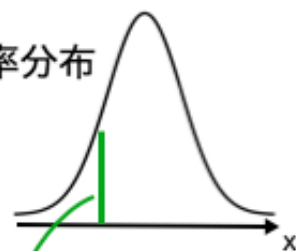
共役事前分布

母数が規定する確率分布



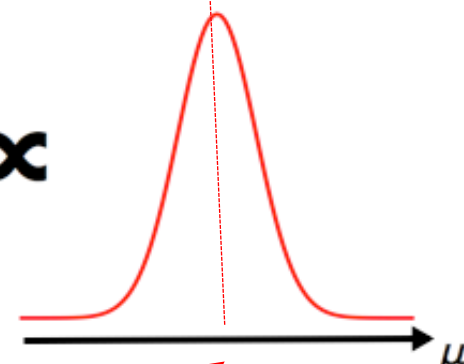
×

尤度



∝

事後分布



同じ形の分布

$$P(w)$$

$$P(w | y, x)$$

データから推定

$w^* = \arg \max_w P(w | y, x)$ が最適なパラメータ.(MAP 推定)

w_* が求まったので, 未知の入力 x_* に対する出力 y_* がどうなるか予測する.

$$\begin{aligned} P(w_* | y_*, x_*) &= \frac{P(w_*)P(y_* | x_*, w_*)}{P(y_* | x_*)} \\ P(y_* | x_*) &= \frac{P(w_*)P(y_* | x_*, w_*)}{P(w_* | y_*, x_*)} \\ &\propto \frac{P(y_* | x_*, w_*)}{P(w_* | y_*, x_*)} \\ &= \frac{N(y_* | x_*, w_*)}{N(w_* | \hat{\mu}|_{y_n=y_*, x_n=x_*}, \hat{\Sigma}^{-1}|_{x_n=x_*})} \\ &= \dots = N(y_* | \mu_*, (\sigma_*^2)^{-1}) \end{aligned}$$

$$\text{ただし, } \mu_* = \mu^T x_*, (\sigma_*^2)^{-1} = (\sigma^2)^{-1} + x_*^T \Sigma^{-1} x_*$$

線形回帰モデルの比較

線形回帰も奥が深い!!

- 関数モデル

$$\mathbf{w}_* = \begin{bmatrix} \bar{y} - \frac{\sigma_{xy}}{\sigma_x^2} \bar{x} \\ \frac{\sigma_{xy}}{\sigma_x^2} \end{bmatrix}$$

$$f(\mathbf{x}) = \mathbf{w}_*^T \mathbf{x}$$

- 識別モデル

$$\mathbf{w}_* = \arg \max_{\mathbf{w}} N(\mathbf{w} \mid \hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}^{-1})$$

$$\hat{\boldsymbol{\Sigma}} = (\sigma)^2 \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T + \boldsymbol{\Sigma}, \quad \hat{\boldsymbol{\mu}} = \hat{\boldsymbol{\Sigma}}^{-1} \{ (\sigma)^2 \sum_{i=1}^N y_i \mathbf{x}_i + \boldsymbol{\Sigma} \boldsymbol{\mu} \}$$

$$f(\mathbf{x}) = N(y_* \mid \mu_*, (\sigma_*^2)^{-1})$$

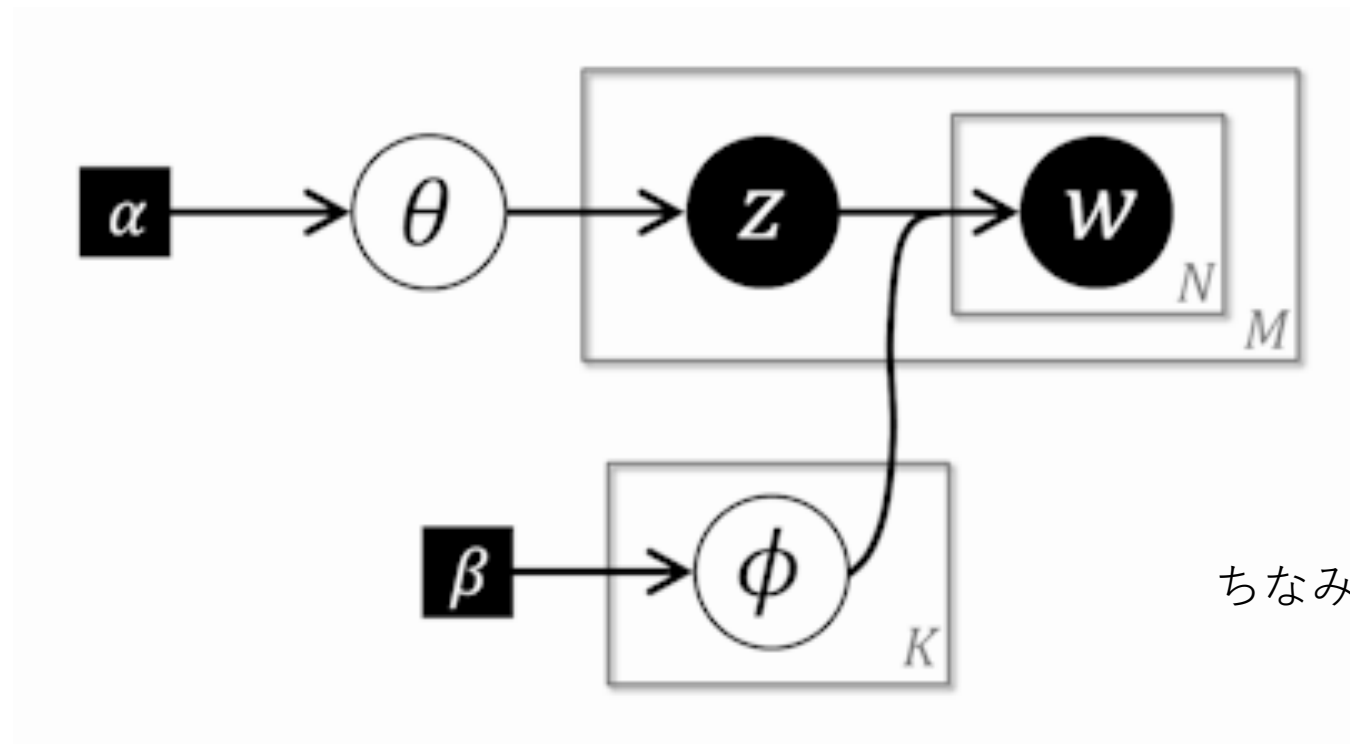
$$\mu_* = \boldsymbol{\mu}^T \mathbf{x}_*, \quad (\sigma_*^2)^{-1} = (\sigma^2)^{-1} + \mathbf{x}_*^T \boldsymbol{\Sigma}^{-1} \mathbf{x}_*$$

- 一般に確率的なアプローチを行うと数学が難しい.
- 特に生成モデルは難しい.
- 一般に事前分布から事後分布を求める方法：
 1. 共役事前分布を使う.
 2. MCMC（マルコフ連鎖モンテカルロ法）を使う.
 3. 変分近似を行う.

2, 3の方法は近似的に確率分布を求める手法

グラフィカルモデル（おまけ）

- グラフによって，確率変数間の条件付き依存構造を示す．



ちなみにDAGになっている．

最後に

線形代数, 解析学, 基礎統計学 は最低勉強しましょう.

ちゃんと理論をやりたい人

集合位相, 複素解析, 測度論, ルベーク積分, 数値解析
関数解析, 位相幾何学, 微分幾何学.