

# 知能機械と自然言語処理

## 知能機械部 第5回

ソフトウェア情報学部

Goutam Chakraborty

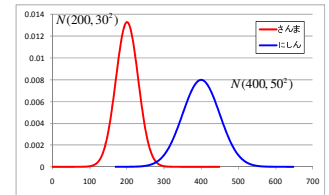
1

## 正規分布の復習

さんまとしんのニラスのデータがあります。特色は魚の重さ(w)です。さんまはクラス0、にしんはクラス1とします。さんまの重さの平均は200gで標準偏差は30です。にしんの平均重さは400gで標準偏差は50gです。

- 0:さんま3000匹
- 1:にしん2000匹

$$P(C_0) = \frac{3000}{5000}, P(C_1) = \frac{2000}{5000}$$



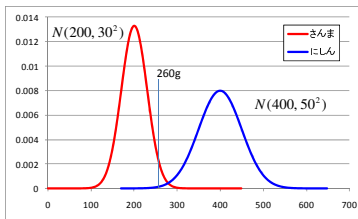
もし特色の重さが260gであれば、その条件付き確率は

$$P(260 \leq w < 260 + \Delta w | C_i)$$

2

## 正規分布復習

$$P(C_0) = \frac{3000}{5000}, P(C_1) = \frac{2000}{5000}$$



if  $w = 260$ あたりの確率を  
 $P(260 \leq w < 260 + \Delta w | C_i)$ とする

3

## 正規分布の確率密度関数から確率を求める簡単な手法

$N(\mu = 200, \sigma^2 = 30^2)$  のとき、

$$f(x) = \frac{1}{30\sqrt{2\pi}} \exp\left(-\frac{1}{2} \times \left(\frac{x-200}{30}\right)^2\right)$$

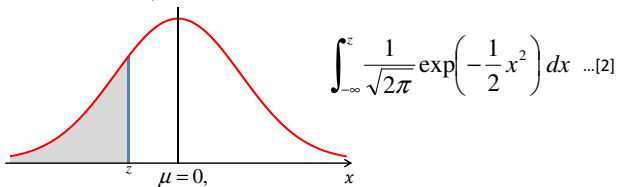
$$P(260 \leq x < (260 + \Delta w) | C_0) = \int_{260}^{260+\Delta w} f(x) dx$$

$$= \int_{260}^{260+\Delta w} \frac{1}{30\sqrt{2\pi}} \exp\left(-\frac{1}{2} \times \left(\frac{x-200}{30}\right)^2\right) dx \dots\dots[1]$$

4

## 標準正規分布の確率密度関数表を参考にして確率を求める方法

前のスライドの式[1]の積分を求めるのは難しい。  
 しかし標準正規分布の確率密度関数表を参考にして確率を簡単に求められる。  
 標準正規分布の  $\mu = 0, \sigma = 1$



$$\int_{-\infty}^z \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} x^2\right) dx \dots[2]$$

上記の図は標準分布の確率密度関数です。網掛けの面積はxの値 $-\infty$ からzまでの確率です。その値は[2]の式で求められます。**標準確率密度関数表**には全てのzの値に対応した網掛けの面積(確率)の値が載っています。

5

## 一般の正規分布から標準正規分布に変更方法

データ  $x$  が、平均値  $\mu$ 、標準偏差  $\sigma$  の一般正規分布データであるとき、  
 $z = (x - \mu) / \sigma$  と加工すると、データは標準正規分布のデータになる。  
 $z$  の平均は0、標準偏差は1になります。

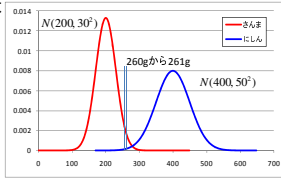
<http://www.six-sigma-material.com/Normal-Distribution.html>

6

### 標準正規分布表使用法

#### C<sub>0</sub>(さんま)の場合赤い分布

平均200、標準偏差30の場合、正規分布の書き方は  
 $x$ から $z$ を求める式は  $Z = (x - 200)/30$   
 $x = 260$ の時  $z = (260 - 200)/30$ です。  
 260g辺り、260gから261g( $\Delta w = 1g$ )の間の重さの  
 さんまの確率は下記通り求められます:



$$P\left(Z \leq \frac{261-200}{30}\right) = P\left(Z \leq \frac{61}{30}\right) \\ = P(Z \leq 2.03) = 0.9788$$

さんまの重さ271gの以下の確率は0.9788です。

$$P\left(Z \leq \frac{260-200}{30}\right) = P\left(Z \leq \frac{60}{30}\right) \\ = P(Z \leq 2.00) = 0.9772$$

さんまの重さ260gの以下の確率は0.9772です。

よって、  
 $\therefore P(C_0 | 260 < w \leq 261) \\ = 0.0016 \times \frac{1}{3} = 0.001$   
 さんまの事前確率 $P(C_0) = 3/5$ です。よって、事後確率  
 $\therefore P(C_0 | 260 < w \leq 261) \\ = 0.0016 \times \frac{3}{5} = 0.001$

同様に、 $C_1$ についても求めると、

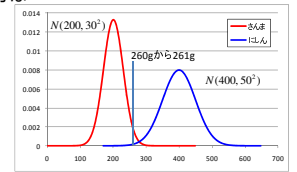
7

### 標準正規分布表使用法

#### C<sub>1</sub>(にしん)の場合青い分布

平均400、標準偏差50の場合、正規分布の書き方は

$x = 400$ の時  $z = (260 - 400)/50$ です。  
 260g辺り、260gから261gの間の重さの  
 にしんの確率は下記通り求められます:



$$P\left(Z \leq \frac{261-400}{50}\right) = P\left(Z \leq \frac{-139}{50}\right) \\ = P(Z \leq -2.78) = 0.0027$$

にしんの重さ261gの以下の確率は0.0027です。

$$P\left(Z \leq \frac{260-400}{50}\right) = P\left(Z \leq \frac{-140}{50}\right) \\ = P(Z \leq -2.8) = 0.0026$$

にしんの重さ260gの以下の確率は0.0026です。

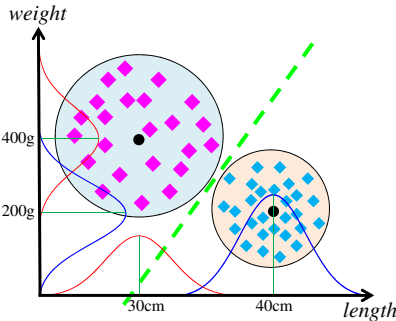
よって、  
 $P(260 < x \leq 261 | C_1) \\ = (0.0027 - 0.0026) = 0.0001$   
 にしんの事前確率 $P(C_1) = 2/5$ です。よって、事後確率  
 $\therefore P(C_1 | 260 < w \leq 261) \\ = 0.0001 \times \frac{5}{2} = 0.00025$

よって、クラス $C_0$ さんまの確率は高い。

$\therefore$ この場合、 $C_0$ に分類される

8

## 2次元の正規分布



各特色を2次元の空間に表示させると、  
 重ならずに分けることができる

9

## d次元の正規分布 (特色の数dの場合)

一次元(特色は一つ)の正規分布

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

#### d次元の正規分布

$$f(\vec{x}) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\vec{x}-\vec{\mu})^T \Sigma^{-1} (\vec{x}-\vec{\mu})\right)$$

$d$ : 次元数(dimension) = 特色の数

$\Sigma$ : 相関行列( $d \times d$ )

10

## d次元の場合ベイズの定理(公式)

$$P(C_0 | 200 \leq \text{weight} < 220 \wedge 35 \leq \text{length} < 37) \\ = P(C_0) \times P(200 \leq \text{weight} < 220 \wedge 35 \leq \text{length} < 37 | C_0)$$

- $d$ 次元の特徴量を表すときは $x_d^n$ で表す  
 $n$ : データの番号,  $d$ : 次元の番号

$$\begin{matrix} x_1^1 = 46.1 & x_2^1 = 180 & \cdots & x_d^1 \\ x_1^2 = 31.6 & x_2^2 = 241 & \cdots & x_d^2 \\ \vdots & & & \vdots \\ x_1^n & x_2^n & \cdots & x_d^n \end{matrix}$$

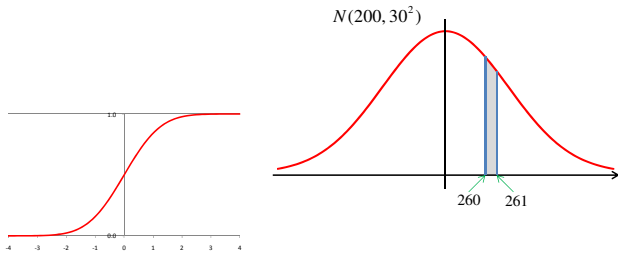
11

STANDARD NORMAL DISTRIBUTION: Table Values Represent AREA to the LEFT of the Z score.

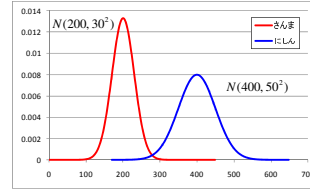
Z	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
-3.9	0.0005	0.0005	0.0004	0.0004	0.0004	0.0004	0.0004	0.0004	0.0003	0.0003
-3.8	0.0007	0.0007	0.0007	0.0006	0.0006	0.0006	0.0006	0.0005	0.0005	0.0005
-3.7	0.0011	0.0010	0.0010	0.0010	0.0009	0.0009	0.0009	0.0008	0.0008	0.0008
-3.6	0.0016	0.0015	0.0015	0.0014	0.0014	0.0013	0.0013	0.0012	0.0012	0.0011
-3.5	0.0023	0.0022	0.0022	0.0021	0.0020	0.0019	0.0019	0.0018	0.0017	0.0017
-3.4	0.0034	0.0032	0.0031	0.0030	0.0029	0.0028	0.0027	0.0026	0.0025	0.0024
-3.3	0.0048	0.0047	0.0045	0.0043	0.0042	0.0040	0.0039	0.0038	0.0036	0.0035
-3.2	0.0069	0.0066	0.0064	0.0062	0.0060	0.0058	0.0056	0.0054	0.0052	0.0050
-3.1	0.0097	0.0094	0.0090	0.0087	0.0084	0.0082	0.0079	0.0076	0.0074	0.0071
-3.0	0.0135	0.0131	0.0126	0.0122	0.0118	0.0114	0.0111	0.0107	0.0104	0.0100
-2.9	0.0187	0.0181	0.0175	0.0169	0.0164	0.0159	0.0154	0.0149	0.0144	0.0139
-2.8	0.0256	0.0248	0.0240	0.0233	0.0226	0.0219	0.0212	0.0205	0.0199	0.0193
-2.7	0.0347	0.0336	0.0326	0.0317	0.0307	0.0298	0.0289	0.0280	0.0272	0.0264
-2.6	0.0466	0.0453	0.0440	0.0427	0.0415	0.0402	0.0391	0.0379	0.0368	0.0357
-2.5	0.0621	0.0604	0.0587	0.0570	0.0554	0.0539	0.0523	0.0508	0.0494	0.0480
-2.4	0.0820	0.0798	0.0776	0.0755	0.0734	0.0714	0.0695	0.0676	0.0657	0.0639
-2.3	0.1072	0.1044	0.1017	0.0990	0.0964	0.0939	0.0914	0.0889	0.0866	0.0842
-2.2	0.1390	0.1355	0.1321	0.1287	0.1255	0.1223	0.1191	0.1160	0.1130	0.1101
-2.1	0.1786	0.1743	0.1700	0.1659	0.1618	0.1578	0.1539	0.1500	0.1463	0.1426
-2.0	0.2275	0.2222	0.2169	0.2118	0.2068	0.2018	0.1970	0.1923	0.1876	0.1831
-1.9	0.2872	0.2807	0.2743	0.2680	0.2619	0.2559	0.2500	0.2442	0.2385	0.2330
-1.8	0.3593	0.3515	0.3438	0.3362	0.3288	0.3216	0.3144	0.3074	0.3005	0.2938
-1.7	0.4457	0.4363	0.4272	0.4182	0.4093	0.4006	0.3920	0.3836	0.3754	0.3673
-1.6	0.5480	0.5370	0.5262	0.5155	0.5050	0.4947	0.4846	0.4746	0.4648	0.4551
-1.5	0.6681	0.6552	0.6426	0.6301	0.6178	0.6057	0.5938	0.5821	0.5705	0.5592
-1.4	0.8076	0.7927	0.7780	0.7636	0.7493	0.7353	0.7215	0.7078	0.6944	0.6811
-1.3	0.9680	0.9510	0.9342	0.9176	0.9012	0.8851	0.8691	0.8534	0.8379	0.8226
-1.2	1.1907	1.1314	1.1123	1.0935	1.0749	1.0565	1.0383	1.0204	1.0027	0.9853
-1.1	1.3987	1.3330	1.3130	1.2924	1.2714	1.2507	1.2302	1.2100	1.1900	1.1702
-1.0	1.5866	1.5625	1.5386	1.5151	1.4917	1.4686	1.4457	1.4231	1.4007	1.3786
-0.9	1.8406	1.8141	1.7879	1.7619	1.7361	1.7106	1.6853	1.6602	1.6354	1.6109
-0.8	2.1186	2.0897	2.0611	2.0327	2.0045	1.9766	1.9489	1.9215	1.8943	1.8673
-0.7	2.4196	2.3885	2.3576	2.3270	2.2965	2.2663	2.2363	2.2065	2.1770	2.1476
-0.6	2.7425	2.7091	2.6761	2.6435	2.6110	2.5785	2.5463	2.5143	2.4825	2.4510
-0.5	3.0854	3.0501	3.0153	2.9806	2.9460	2.9116	2.8774	2.8434	2.8096	2.7760
-0.4	3.4438	3.4090	3.3724	3.3360	3.2997	3.2636	3.2276	3.1918	3.1561	3.1207
-0.3	3.8209	3.7828	3.7448	3.7070	3.6693	3.6317	3.5942	3.5569	3.5197	3.4827
-0.2	4.2074	4.1683	4.1294	4.0905	4.0517	4.0129	3.9743	3.9358	3.8974	3.8591
-0.1	4.6017	4.5620	4.5224	4.4828	4.4433	4.4038	4.3644	4.3251	4.2858	4.2465
0.0	5.0000	4.9601	4.9202	4.8803	4.8405	4.8006	4.7608	4.7210	4.6812	4.6414

12

ちなみに、ExcelではNORMSDISTという関数で  
上の式の計算ができる



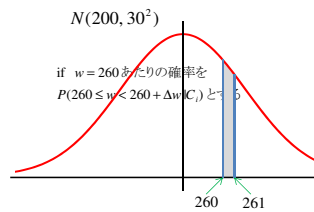
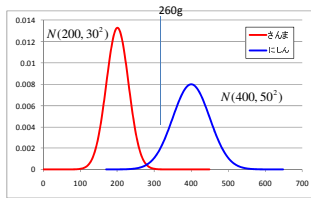
13



14

## 正規分布復習

$$P(C_0) = \frac{3000}{5000}, \quad P(C_1) = \frac{2000}{5000}$$



if  $w = 260$ あたりの確率を  
 $P(260 \leq w < 260 + \Delta w | C_1)$  とする

15