
Lab 8 (solution)

Probabilistic Reasoning - Gaussian distribution and Bayes Theorem

FTP MachLe MSE
HS 2023

Machine Learning
WÜRC

The central *paradigm of probabilistic reasoning* is to identify all relevant variables x_1, \dots, x_N in the environment, and make a probabilistic model $p(x_1, \dots, x_N)$ of their interaction. Reasoning (*inference*) is then performed by introducing *evidence* that sets variables in known states, and subsequently computing probabilities of interest, *conditioned on this evidence*. The rules of probability, combined with Bayes' rule make a reasoning system complete.

After this unit, ...

Lernziele/Kompetenzen

- you have repeated the *basic rules of probability theory*.
 - you know the difference between a *joint* and a *conditional probability* distribution.
 - you know how to apply *Bayes Theorem* to calculate the *posterior* probability distribution for simple discrete examples. You can name the *prior* probability distribution, the *likelihood function*, the *evidence*, and you know how to *marginalize* over a joint probability distribution.
 - you know the basic properties of a *multivariate Gaussian* probability distribution. You can plot a 2D Gaussian probability distribution given the *mean vector* μ and the covariance matrix Σ .
 - you can *sample* data points from a given multivariate gaussian distribution.
 - you can explain the *naïve Bayes classifier* to your classmates and to your teacher.
-

1. Supervised Bayesian Learning [M,II]

The table below contains the result of a market survey for a promotion for different items such as a magazine, a watch and a life insurance and a credit card insurance. 10 people were interviewed and asked whether they would buy such items.

Use this count table for supervised Bayesian learning. The output attribute is *sex* with possible values **male** and **female**. Consider an individual who has said *no* to the life insurance promotion, *yes* to the magazine promotion, *yes* to the watch promotion and *yes* to the credit card insurance. Use the values in the table together with the Naive Bayes classifier to determine which of a,b,c or d represents the probability that this individual is male. $p(E)$ is the marginal distribution.

	Magazine		Watch		Life Insurance		Credit Card	
	male	female	male	female	male	female	male	female
yes	4	3	2	2	2	3	2	1
no	2	1	4	2	4	1	4	3

Welche der folgenden Aussagen sind wahr und welche falsch?	wahr	falsch
a) $p(\text{sex} = \text{male} \dots) = \frac{4}{6} \cdot \frac{2}{6} \cdot \frac{2}{6} \cdot \frac{2}{6} \cdot \frac{6}{10} \cdot \frac{1}{p(E)}$	<input type="radio"/>	<input checked="" type="radio"/>
b) $p(\text{sex} = \text{male} \dots) = \frac{4}{6} \cdot \frac{2}{6} \cdot \frac{3}{4} \cdot \frac{2}{6} \cdot \frac{3}{4} \cdot \frac{1}{p(E)}$	<input type="radio"/>	<input checked="" type="radio"/>
c) $p(\text{sex} = \text{male} \dots) = \frac{4}{6} \cdot \frac{4}{6} \cdot \frac{2}{6} \cdot \frac{2}{6} \cdot \frac{6}{10} \cdot \frac{1}{p(E)}$	<input checked="" type="radio"/>	<input type="radio"/>
d) $p(\text{sex} = \text{male} \dots) = \frac{2}{6} \cdot \frac{4}{6} \cdot \frac{4}{6} \cdot \frac{2}{6} \cdot \frac{4}{10} \cdot \frac{1}{p(E)}$	<input type="radio"/>	<input checked="" type="radio"/>

2. Hamburger and Bayes Rule [A,II]

- a) The probability that a hamburger eater HE will have Kreuzfeld-Jacob disease given the prior $p(\text{KJ})$ and the marginal $p(\text{HE}) = \sum_{\text{KJ}} p(\text{HE}|\text{KJ}) \cdot p(\text{KJ})$ is:

$$p(\text{KJ}|\text{HE}) = \frac{p(\text{HE}, \text{KJ})}{p(\text{HE})} = \frac{p(\text{HE}|\text{KJ}) \cdot p(\text{KJ})}{p(\text{HE})} \quad (1)$$

$$= \frac{p(\text{HE}|\text{KJ}) \cdot p(\text{KJ})}{\sum_{\text{KJ}} p(\text{HE}|\text{KJ}) \cdot p(\text{KJ})} = \frac{\frac{9}{10} \cdot \frac{1}{100000}}{\frac{1}{2}} \approx \underline{\underline{1.8 \cdot 10^{-5}}} \quad (2)$$

- b) If the fraction of people eating hamburgers was rather small, $p(\text{HE}) = 0,001$, what is the probability that a regular hamburger eater will have Kreuzfeld-Jacob disease?

$$p(\text{KJ}|\text{HE}) = \frac{p(\text{HE}, \text{KJ})}{p(\text{HE})} = \frac{p(\text{HE}|\text{KJ}) \cdot p(\text{KJ})}{p(\text{HE})} \quad (3)$$

$$= \frac{\frac{9}{10} \cdot \frac{1}{100000}}{\frac{1}{1000}} = \underline{\underline{9 \cdot 10^{-3}}} \approx 1\% \quad (4)$$

3. Naïve Bayes Classifier [A,II]

a) Python Code: Naive Bayes prior and likelihoods

```
p_y = 4.0/10; # p(y) = 4/10
# p(xi=1 | y=-1)
p_x1_y0 = 3.0/6;
p_x2_y0 = 5.0/6;
p_x3_y0 = 4.0/6;
p_x4_y0 = 5.0/6;
p_x5_y0 = 2.0/6;

# p(xi=1 | y=+1)
p_x1_y1 = 3.0/4;
p_x2_y1 = 0.0/4;
p_x3_y1 = 3.0/4;
p_x4_y1 = 2.0/4;
p_x5_y1 = 1.0/4;
```

b) Python Code: Naive Bayes classification decisions

```
f_y1_00000 = p_y*(1-p_x1_y1)*(1-p_x2_y1)*(1-p_x3_y1)*
(1-p_x4_y1)*(1-p_x5_y1)
print("f_y1_00000 = ",f_y1_00000)

f_y0_00000 = (1-p_y)*(1-p_x1_y0)*(1-p_x2_y0)*(1-p_x3_y0)*
(1-p_x4_y0)*(1-p_x5_y0)
print("f_y0_00000 = " , f_y0_00000)

if (f_y1_00000 > f_y0_00000):
    print("Predict class +1")
else:
    print ("Predict class -1")
print("\n")

f_y1_11010 = p_y*(p_x1_y1)*(p_x2_y1)*(1-p_x3_y1)*
(p_x4_y1)*(1-p_x5_y1)
print ("f_y1_11010 = ",f_y1_11010)

f_y0_11010 = (1-p_y)*(p_x1_y0)*(p_x2_y0)*(1-p_x3_y0)*
(p_x4_y0)*(1-p_x5_y0)
print("f_y0_11010 = ",f_y0_11010)

if (f_y1_11010 > f_y0_11010):
    print("Predict class +1")
else:
    print ("Predict class -1")
```

The numerical solution $x = \{00000\}$ is:

$$\begin{aligned}f(y = +1|00000) &= 0.009375000000000001 \\f(y = -1|00000) &= 0.0018518518518518515 \\f(y = +1|00000) &> f(y = -1|00000) \implies \hat{y} = +1\end{aligned}$$

The numerical solution for $x = \{11010\}$ is:

$$\begin{aligned} f(y = +1|11010) &= 0.0 \\ f(y = -1|11010) &= 0.046296296296296315 \\ f(y = +1|11010) &< f(y = -1|11010) \implies \hat{y} = -1 \end{aligned}$$

c) Python Code: Naive Bayes posterior probabilities

```
# p(y=1|00000) =
print ("p(y=1|00000) =", f_y1_00000 / (f_y1_00000 + f_y0_00000))

# p(y=1|11010) =
print ("p(y=1|11010) =", f_y1_11010 / (f_y1_11010 + f_y0_11010))
```

The Naive Bayes posterior probabilities are:

$$\begin{aligned} p(y = 1|00000) &= 0.8350515463917526 \\ p(y = 1|11010) &= 0.0 \end{aligned}$$

- d)** A *Bayes classifier using a joint distribution model* for $p(x_1, \dots, x_5|y = c)$ would have $2^5 - 1 = 31$ degrees of freedom (independent probabilities) to estimate. But here, we have only 6 data points from class $y = -1$, and 4 data points from class $y = +1$. Thus these models would assign zero probability to many feature combinations, and *would probably not generalize well* to new data.
- e)** No, we do not need to re-train the model by estimating new probabilities. Due to the conditional independence assumptions of naïve Bayes, it is optimal to simply ignore $p(x_1|y)$, and use the previously estimated probabilities of the other four features when computing $p(y|x_2, x_3, x_4, x_5)$. If you did recompute $p(x_i|y)$ using the formulas above, it is easy to verify that the values would not change.

4. Passenger Scanner [A,II]

The detector is such that 95% of all terrorists are identified as terrorists

$$p(\text{label} = \text{true} \mid \text{terr} = \text{true}) = 0.95 \quad (5)$$

We can infer from the former that only 5% of terrorists got mislabeled as good people (false negative):

$$p(\text{label} = \text{false} \mid \text{terr} = \text{true}) = 0.05 \quad (6)$$

Furthermore 95% of all upstanding citizens are identified as such (non terrorists). Therefore, the false positive rate is:

$$p(\text{label} = \text{true} \mid \text{terr} = \text{false}) = 0.05 \quad (7)$$

Assuming that the informant is correct:

$$p(\text{terr} = \text{true}) = \frac{1}{100} \quad (8)$$

We can then infer that $\frac{99}{100}$ are not terrorists:

$$p(\text{terr} = \text{false}) = \frac{99}{100} \quad (9)$$

We can find out the probability that the person picked is a terrorist using *Bayes' Rule*:

$$p(\text{terr} = \text{true} \mid \text{label} = \text{true}) = \frac{p(\text{label} = \text{true} \mid \text{terr} = \text{true})p(\text{terr} = \text{true})}{p(\text{label} = \text{true})} \quad (10)$$

To find out the denominator, we can *marginalize* over the joint distribution of $p(\text{label} = \text{true})$ and $p(\text{terr})$:

$$\begin{aligned} p(\text{label} = \text{true}) &= \sum_t p(\text{label} = \text{true}, \text{terr} = t) \\ &= \sum_t p(\text{label} = \text{true} \mid \text{terr} = t) \cdot p(\text{terr} = t) \end{aligned}$$

We make the summation explicit:

$$\begin{aligned} p(\text{label} = \text{true}) &= p(\text{label} = \text{true} \mid \text{terr} = \text{true})p(\text{terr} = \text{true}) + \\ &\quad p(\text{label} = \text{true} \mid \text{terr} = \text{false})p(\text{terr} = \text{false}) \\ &= 0.95 \cdot 0.01 + 0.05 \cdot 0.99 = \underline{\underline{0.059}} \end{aligned}$$

Substituting all the values into Bayes' Rule, we find out that even though the detector is 'reliable' we still get a low posterior probability that the person suspected of being a terrorist is actually a terrorist:

$$\begin{aligned} p(\text{terr} = \text{true} \mid \text{label} = \text{true}) &= \frac{p(\text{label} = \text{true} \mid \text{terr} = \text{true})p(\text{terr} = \text{true})}{p(\text{label} = \text{true})} \\ &= \frac{0.95(0.01)}{0.059} = \underline{\underline{0.1610169492}} \end{aligned}$$

5. Bivariate Gaussian Distribution [A,II]

a) The *bivariate normal density* $p(\mathbf{x}) = p(x_a, x_b)$ is defined by

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2} \sqrt{\det(\boldsymbol{\Sigma})}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

In this equation, we have $D = 2$ and

$$\begin{aligned} \mathbf{x} &= \begin{pmatrix} x_a \\ x_b \end{pmatrix} \\ \boldsymbol{\mu} &= \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} = \begin{pmatrix} 0 \\ 2 \end{pmatrix} \\ \boldsymbol{\Sigma} &= \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} = \begin{pmatrix} 2 & \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} & 1 \end{pmatrix} \end{aligned}$$

$$|\Sigma| = \det(\Sigma) = 2 \cdot 1 - \left(\frac{\sqrt{2}}{2}\right)^2 = \frac{3}{2}$$

$$\sqrt{|\Sigma|} = \sqrt{\frac{3}{2}}$$

The inverse matrix of a 2×2 -matrix can be calculated by:

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

$$A^{-1} = \frac{1}{\det(A)} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$$

In this case, we get the following result for the precision matrix Λ :

$$\Lambda = \frac{2}{3} \begin{pmatrix} 1 & -\frac{\sqrt{2}}{2} \\ -\frac{\sqrt{2}}{2} & 2 \end{pmatrix}$$

- b)** The squared generalized distance expression, i.e. the *Mahalanobis distance* can be written as:

$$\begin{aligned} \Delta &= (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \\ &= \begin{pmatrix} x_a \\ x_b - 2 \end{pmatrix}^T \frac{2}{3} \begin{pmatrix} 1 & -\frac{\sqrt{2}}{2} \\ -\frac{\sqrt{2}}{2} & 2 \end{pmatrix} \begin{pmatrix} x_a \\ x_b - 2 \end{pmatrix} \\ &= \frac{2}{3} (x_a^2 - \sqrt{2}x_a(x_b - 2) + 2(x_b - 2)^2) \end{aligned}$$

as a function of x_a and x_b .

- c)** The joint probability density is explicitly given by:

$$\begin{aligned} p(\mathbf{x}) = p(x_a, x_b) &= \frac{1}{(2\pi)^{2/2} \sqrt{3/2}} \cdot \exp \left\{ -\frac{1}{2} \begin{pmatrix} x_a \\ x_b - 2 \end{pmatrix}^T \frac{2}{3} \begin{pmatrix} 1 & -\frac{\sqrt{2}}{2} \\ -\frac{\sqrt{2}}{2} & 2 \end{pmatrix} \begin{pmatrix} x_a \\ x_b - 2 \end{pmatrix} \right\} \\ &= \frac{1}{\sqrt{6\pi}} \cdot \exp \left\{ -\frac{1}{3} (x_a^2 - \sqrt{2}x_a(x_b - 2) + 2(x_b - 2)^2) \right\} \end{aligned}$$

- d)** We calculate the *eigenvalues* $\lambda_{1,2}$ and the *eigenvectors* $\mathbf{u}_{1,2}$ of the covariance matrix Σ using `np.linalg.eig`.

```
import numpy as np
w, v = np.linalg.eig(np.array([[2, np.sqrt(2)/2], [np.sqrt(2)/2, 1]]))

print(w)
print(v)

[2.3660254  0.6339746]
[[ 0.88807383 -0.45970084]
 [ 0.45970084  0.88807383]]
```

The eigenvalues of Σ are $(\lambda_1, \lambda_2) = (2.3660254, 0.6339746)$ with eigenvectors:

$$(\mathbf{u}_1, \mathbf{u}_1) = \begin{pmatrix} 0.88807383 & -0.45970084 \\ 0.45970084 & 0.88807383 \end{pmatrix}$$

e) The Python code could look like this:

```
# -*- coding: utf-8 -*-
"""
Created on Tue Jan 21 19:52:21 2020
@author: wuersch
"""

import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import multivariate_normal

# Mean vector and covariance matrix
mu      = np.array([0., 2.])
Sigma   = np.array([[2, np.sqrt(2)/2], [np.sqrt(2)/2, 1]])

F = multivariate_normal(mu, Sigma)

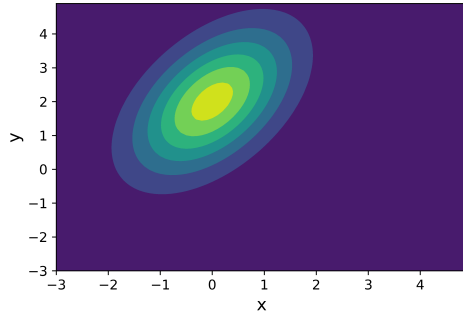
#draw random samples from the multivariate distribution
#and try to reconstruct the gaussian distribution

NSamples=10000
x, y = np.mgrid[-3:5:.1, -3:5:.1]

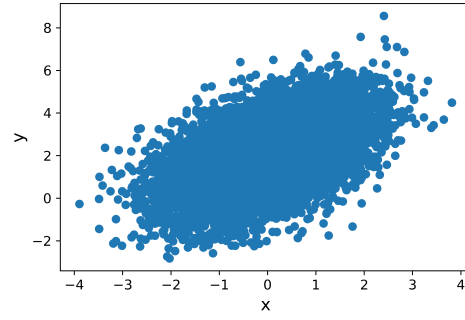
pos      = np.dstack((x, y))
MVGauss  = multivariate_normal(mu, Sigma)
MVGSamples = MVGauss.rvs(size=NSamples)
XS        = MVGSamples[:,0]
YS        = MVGSamples[:,1]

fig2 = plt.figure('Using Scikit Learn')
ax2  = fig2.add_subplot(111)
ax2.contourf(x, y, F.pdf(pos))
ax2.set_xlabel("x")
ax2.set_ylabel("y")
plt.savefig('Gauss3D_2.png', dpi=600)

fig3 = plt.figure('Using Scikit Learn to draw random samples')
ax3  = fig3.add_subplot(111)
ax3.scatter(XS, YS)
ax3.set_xlabel("x")
ax3.set_ylabel("y")
plt.savefig('Gauss3D_3.png', dpi=600)
```



(a) bivariate gaussian



(b) 10'000 samples

- f) We calculate the *conditional probability* $p(x_a|x_b)$. In general, the conditional of a multi-variate Gaussian distribution is again gaussian with a mean $\mu_{a|b}$ and a covariance matrix $\Sigma_{a|b} = \Lambda_{aa}^{-1}$

$$\Lambda = \Sigma^{-1} = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$

$$\begin{aligned} p(\mathbf{x}_a|\mathbf{x}_b) &= \mathcal{N}(\mathbf{x}_a|\mu_{a|b}, \Lambda_{aa}^{-1}) \\ \mu_{a|b} &= \mu_a - \Lambda_{aa}^{-1} \Lambda_{ab}(\mathbf{x}_b - \mu_b) \\ \Sigma_{a|b} &= \Lambda_{aa}^{-1} \end{aligned}$$

$$\begin{aligned} \mu_{a|b} &= \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (\mathbf{x}_b - \mu_b) \\ &= 0 + \frac{\sqrt{2}}{2} \cdot (1)^{-1} (x_b - 2) \\ &= \frac{\sqrt{2}}{2} (x_b - 2) \end{aligned}$$

$$\begin{aligned} \Sigma_{a|b} &= \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba} \\ &= 2 - \frac{\sqrt{2}}{2} \cdot (1)^{-1} \cdot \frac{\sqrt{2}}{2} \\ &= 2 - \frac{1}{2} = \frac{3}{2} \\ &= \Lambda_{aa}^{-1} = \left(\frac{2}{3}\right)^{-1} = \frac{3}{2} \end{aligned}$$

$$\begin{aligned} p(\mathbf{x}_a|\mathbf{x}_b) &= \mathcal{N}(\mathbf{x}_a|\mu_{a|b}, \Lambda_{aa}^{-1}) \\ &= \mathcal{N}\left(\mathbf{x}_a \middle| \frac{\sqrt{2}}{2} (x_b - 2), \frac{3}{2}\right) \end{aligned}$$

6. Weather in London [A,II]

- a) Assuming that the prior probability it rained yesterday is 0.5, what is the probability that it was raining yesterday given that it's sunny today?

$$p(\text{yesterday} = \text{rain}) = 50\%$$

I infer from this that the probability of being sunny yesterday is also 50%:

$$p(\text{yesterday} = \text{sun}) = 50\%$$

$$\begin{aligned} p(\text{yesterday} = \text{rain} \mid \text{today} = \text{sun}) \\ = \frac{p(\text{today} = \text{sun} \mid \text{yesterday} = \text{rain})p(\text{yesterday} = \text{rain})}{p(\text{today} = \text{sun})} \end{aligned}$$

$$\begin{aligned} p(\text{today} = \text{sun}) &= \sum_y p(\text{today} = \text{sun}, \text{yesterday} = y) \\ &= \sum_y p(\text{today} = \text{sun} \mid \text{yesterday} = y)p(\text{yesterday} = y) \\ &= p(\text{today} = \text{sun} \mid \text{yesterday} = \text{sun}) \cdot p(\text{yesterday} = \text{sun}) \\ &\quad + p(\text{today} = \text{sun} \mid \text{yesterday} = \text{rain}) \cdot p(\text{yesterday} = \text{rain}) \\ &= 0.40 \cdot 0.5 + 0.30 \cdot 0.50 = 0.20 + 0.15 = \underline{0.35} \end{aligned}$$

Thus, the probability of raining yesterday given that today is sunny:

$$p(\text{yesterday} = \text{rain} \mid \text{today} = \text{sun}) = \frac{0.15}{0.35} = \underline{42.86\%}$$

- b)** If the weather follows the same pattern as above, day after day, what is the probability that it will rain on any day (based on an effectively infinite number of days of observing the weather)?

On any day, not considering whether it rained or not the day before, the probability of raining is:

$$\begin{aligned} p(\text{today} = \text{rain}) &= \sum_y p(\text{today} = \text{rain}, \text{yesterday} = y) \\ &= \sum_y p(\text{today} = \text{rain} \mid \text{yesterday} = y)p(\text{yesterday} = y) \\ &= p(\text{today} = \text{rain} \mid \text{yesterday} = \text{rain})p(\text{yesterday} = \text{rain}) \\ &\quad + p(\text{today} = \text{rain} \mid \text{yesterday} = \text{sun})p(\text{yesterday} = \text{sun}) \\ &= 0.7 \cdot 0.5 + 0.6 \cdot 0.5 = \underline{0.65} \end{aligned}$$

- c)** Use the result from b) above as a new prior probability of rain yesterday and recompute the probability that it was raining yesterday given that it's sunny today.

$$p(\text{yesterday} = \text{rain}) = 0.65$$

Therefore:

$$p(\text{yesterday} = \text{sun}) = 0.35$$

$$p(\text{yesterday} = \text{rain} \mid \text{today} = \text{sun}) = \frac{p(\text{today} = \text{sun} \mid \text{yesterday} = \text{rain})p(\text{yesterday} = \text{rain})}{p(\text{today} = \text{sun})}$$

$$p(\text{today} = \text{sun}) = \sum_y p(\text{today} = \text{sun}, \text{yesterday} = y)$$

$$\begin{aligned}
&= \sum_y p(\text{today} = \text{sun} \mid \text{yesterday} = y) p(\text{yesterday} = y) \\
&= p(\text{today} = \text{sun} \mid \text{yesterday} = \text{sun}) p(\text{yesterday} = \text{sun}) \\
&\quad + p(\text{today} = \text{sun} \mid \text{yesterday} = \text{rain}) p(\text{yesterday} = \text{rain}) \\
&= 0.40 \cdot 0.35 + 0.30 \cdot 0.65 \\
&= 0.14 + 0.195 = \underline{\underline{0.335}}
\end{aligned}$$

$$p(\text{yesterday} = \text{rain} \mid \text{today} = \text{sun}) = \frac{0.30 \cdot 0.65}{0.335} = \underline{\underline{58.21\%}}$$

Very interesting. The probability of raining has increased a little bit after updating the prior probabilities. Since it was likely that it rained yesterday, it's slightly more likely that it will rain today.

7. Inspector Clouseau [A,II]

a) Using b for the two states of B and m for the two states of M ,

$$p(B|K) = \sum_m p(B, m|K) = \sum_m \frac{p(B, m, K)}{p(K)} = \frac{p(B) \sum_m p(K|B, m)p(m)}{\sum_b p(b) \sum_m p(K|b, m)p(m)}$$

Plugging in the values we have

$$\begin{aligned}
p(B = \text{murderer} | \text{knife used}) &= \frac{\frac{6}{10} \left(\frac{2}{10} \cdot \frac{1}{10} + \frac{8}{10} \cdot \frac{6}{10} \right)}{\frac{6}{10} \left(\frac{2}{10} \cdot \frac{1}{10} + \frac{8}{10} \cdot \frac{6}{10} \right) + \frac{4}{10} \left(\frac{2}{10} \cdot \frac{2}{10} + \frac{8}{10} \cdot \frac{3}{10} \right)} \\
&= \frac{300}{412} \approx \underline{\underline{0.73}}
\end{aligned}$$

Remark: The role of $p(\text{knife used})$ in the Inspector Clouseau example can cause some confusion. In the above, $p(\text{knife used})$ is computed to be 0.412.

$$p(\text{knife used}) = \sum_b p(b) \cdot \sum_m p(\text{knife used}|b, m) \cdot p(m) \approx \underline{\underline{0.412}} \quad (11)$$

But surely, $p(\text{knife used}) = 1$, since this is given in the question! Note that the quantity $p(\text{knife used})$ relates to the *prior probability* the model assigns to the knife being used (in the absence of any other information). If we know that the knife is used, then the posterior $p(\text{knife used}) = 1$.

8. Factorization of a multivariate probability distribution [A,II]

A belief network (BN) is a distribution of the form

$$p(x_1, \dots, x_N) = \prod_{i=1}^N p(x_i | \text{pa}(x_i)) \quad (12)$$

where $\text{pa}(x_i)$ represent the *parental* variables of variable x_i . Represented as a directed graph, with an arrow pointing from a parent variable to child variable, a *belief network* corresponds to a *Directed Acyclic Graph* (DAG), with the i^{th} node in the graph corresponding to the factor $p(x_i|\text{pa}(x_i))$. In general, two different graphs may represent the same independence assumptions. If one wishes to make independence assumptions, then the choice of factorisation becomes significant.

The observation that any distribution may be written in the *cascade form*, gives an algorithm for constructing a BN on variables x_1, \dots, x_n : write down the n -node cascade graph; label the nodes with the variables in any order; now each successive independence statement corresponds to deleting one of the edges. More formally, this corresponds to an ordering of the variables which, without loss of generality, we may write as x_1, \dots, x_n . Then, from Bayes' rule, we have

$$p(x_1, \dots, x_N) = p(x_1|x_2, \dots, x_N) \cdot p(x_2, \dots, x_N) \quad (13)$$

$$= p(x_1|x_2, \dots, x_N) \cdot p(x_2|x_3, \dots, x_N) \cdot p(x_3, \dots, x_N) \quad (14)$$

$$= p(x_n) \cdot \prod_{i=1}^{N-1} p(x_i|x_{i+1}, \dots, x_N) \quad (15)$$

The representation of any BN is therefore a Directed Acyclic Graph (DAG). Every probability distribution can be written as a BN, even though it may correspond to a fully connected 'cascade' DAG. The particular role of a BN is that the structure of the DAG corresponds to a *set of conditional independence assumptions*, namely which ancestral parental variables are sufficient to specify each conditional probability table. Note that this does not mean that non-parental variables have no influence.

9. Conditional distribution of a bivariate Gaussian distribution [A,II]

- We start by calculating the precision matrix $\mathbf{\Lambda} = \mathbf{\Sigma}^{-1}$, the inverse of the covariance matrix $\mathbf{\Sigma}$. We use the fact that the inverse of a 2×2 -matrix \mathbf{A} is given by:

$$\mathbf{A} = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \quad (16)$$

$$\mathbf{A}^{-1} = \frac{1}{ad - bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix} \quad (17)$$

In our case, we have

$$\mathbf{\Sigma} = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix} \quad (18)$$

$$\mathbf{\Lambda} = \mathbf{\Sigma}^{-1} = \frac{1}{1 - \rho^2} \begin{pmatrix} \frac{1}{\sigma_1^2} & -\frac{\rho}{\sigma_1\sigma_2} \\ -\frac{\rho}{\sigma_1\sigma_2} & \frac{1}{\sigma_2^2} \end{pmatrix} \quad (19)$$

- Now, we concentrate only on the *quadratic term* \square in the exponential without the factor $-\frac{1}{2}$ which is given by:

$$\square = \frac{1}{1 - \rho^2} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix}^T \begin{pmatrix} \frac{1}{\sigma_1^2} & -\frac{\rho}{\sigma_1\sigma_2} \\ -\frac{\rho}{\sigma_1\sigma_2} & \frac{1}{\sigma_2^2} \end{pmatrix} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix} \quad (20)$$

$$= \frac{1}{1-\rho^2} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix}^T \cdot \begin{pmatrix} \frac{1}{\sigma_1^2}(x_1 - \mu_1) - \frac{\rho}{\sigma_1\sigma_2}(x_2 - \mu_2) \\ \frac{1}{\sigma_2^2}(x_2 - \mu_2) - \frac{\rho}{\sigma_1\sigma_2}(x_1 - \mu_1) \end{pmatrix} \quad (21)$$

$$= \frac{1}{1-\rho^2} \left\{ \frac{1}{\sigma_1^2}(x_1^2 - 2\mu_1x_1 + \mu_1^2) - \frac{\rho}{\sigma_1\sigma_2}(x_1x_2 - x_1\mu_2 - x_2\mu_1 + \mu_1\mu_2) \right\} \quad (22)$$

$$+ \frac{1}{1-\rho^2} \left\{ \frac{1}{\sigma_2^2}(x_2^2 - 2\mu_2x_2 + \mu_2^2) - \frac{\rho}{\sigma_1\sigma_2}(x_1x_2 - x_1\mu_2 - x_2\mu_1 + \mu_1\mu_2) \right\} \quad (23)$$

- Now, we only consider terms quadratic and linear in x_1 and get:

$$\square = \frac{1}{1-\rho^2} \left\{ \frac{x_1^2}{\sigma_1^2} - x_1 \cdot \left(\frac{2\mu_1}{\sigma_1^2} + \frac{2\rho}{\sigma_1\sigma_2}(x_2 - \mu_2) \right) + \dots \right\} \quad (24)$$

$$= \frac{1}{1-\rho^2} \frac{1}{\sigma_1^2} \left\{ x_1^2 - x_1 \cdot \left(2\mu_1 + \frac{2\rho\sigma_1^2}{\sigma_1\sigma_2}(x_2 - \mu_2) \right) + \dots \right\} \quad (25)$$

- By completing the square, we finally get:

$$\square = \frac{1}{1-\rho^2} \frac{1}{\sigma_1^2} \left\{ \left(x_1 - \left(\mu_1 + \frac{\rho\sigma_1}{\sigma_2}(x_2 - \mu_2) \right) \right)^2 + \dots \right\} \quad (26)$$

- By identifying the coefficient before the quadratic term as the inverse of the variance $\tilde{\sigma}$ of the conditional probability $p(x_1|x_2)$ and the shift as the mean $\tilde{\mu}$ of the conditional probability density, we find:

$$\tilde{\mu} = \mu_1 + \rho \frac{\sigma_1}{\sigma_2} \cdot (x_2 - \mu_2) \quad (27)$$

$$\tilde{\sigma} = (1 - \rho^2) \cdot \sigma_1^2 \quad (28)$$