

Introduction to Artificial Intelligence (INFO8006)

Exercise session 4

Maximum likelihood estimation

Given a set of i.i.d. observations $\mathcal{D} = \{x_1, \dots, x_N\}$, a set of unknown parameters $\theta = [\theta_1, \dots, \theta_K]$ and the likelihood function $P(x_i | \theta)$ of one observation given the parameters, we derive the likelihood of the parameters for this set of observations as

$$P(\mathcal{D} | \theta) = \prod_{i=1}^N P(x_i | \theta).$$

From which we can recover the maximum likelihood estimate θ^* of the parameters as

$$\theta^* = \underset{\theta}{\operatorname{argmax}} P(\mathcal{D} | \theta).$$

This can typically be found by cancelling the derivative of the associated log-likelihood w.r.t. each parameter

$$\frac{\partial LL(\mathcal{D}; \theta)}{\partial \theta_k} = 0, \quad \forall k.$$

Bayesian learning and maximum a posteriori

We can treat parameters as random variables to incorporate uncertainty about their values. To do so, we have to specify a prior distribution $P(\theta)$ over the parameters. When new observations \mathcal{D} are collected, the distribution over parameters can be updated, leading to the posterior

$$P(\theta | \mathcal{D}) \propto P(\mathcal{D} | \theta)P(\theta).$$

When the latter is not analytically tractable, we can still compute the maximum a posteriori

$$\theta^* = \underset{\theta}{\operatorname{argmax}} P(\theta | \mathcal{D}).$$

Cheat sheet for Gaussian models

From the joint

$$p\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} = \mathcal{N}\left(\begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix}, \begin{pmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^T & \mathbf{B} \end{pmatrix}\right),$$

the marginal and conditional distributions are given by

$$\begin{aligned} p(\mathbf{x}) &= \mathcal{N}(\mathbf{a}, \mathbf{A}) \\ p(\mathbf{y}) &= \mathcal{N}(\mathbf{b}, \mathbf{B}) \\ p(\mathbf{x} | \mathbf{y}) &= \mathcal{N}(\mathbf{a} + \mathbf{CB}^{-1}(\mathbf{y} - \mathbf{b}), \mathbf{A} - \mathbf{CB}^{-1}\mathbf{C}^T) \\ p(\mathbf{y} | \mathbf{x}) &= \mathcal{N}(\mathbf{b} + \mathbf{C}^T\mathbf{A}^{-1}(\mathbf{x} - \mathbf{a}), \mathbf{B} - \mathbf{C}^T\mathbf{A}^{-1}\mathbf{C}). \end{aligned}$$

On the other hand, from

$$\begin{aligned} p(\mathbf{x}) &= \mathcal{N}(\mathbf{m}, \mathbf{P}) \\ p(\mathbf{y} | \mathbf{x}) &= \mathcal{N}(\mathbf{H}\mathbf{x} + \mathbf{u}, \mathbf{R}), \end{aligned}$$

we can recover the joint distribution as

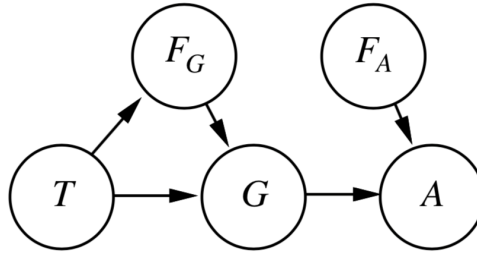
$$p\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} = \mathcal{N}\left(\begin{pmatrix} \mathbf{m} \\ \mathbf{H}\mathbf{m} + \mathbf{u} \end{pmatrix}, \begin{pmatrix} \mathbf{P} & \mathbf{P}\mathbf{H}^T \\ \mathbf{H}\mathbf{P} & \mathbf{H}\mathbf{P}\mathbf{H}^T + \mathbf{R} \end{pmatrix}\right),$$

In session exercises: Ex. 1, Ex. 2

Exercise 1 Nuclear power plant (AIMA, Ex 14.11)

In your local nuclear power station, there is an alarm that senses when a temperature gauge exceeds a given threshold. The gauge measures the temperature of the core. Consider the Boolean variables A (alarm sounds), F_A (alarm is faulty), and F_G (gauge is faulty) and the multi-valued nodes G (gauge reading) and T (actual core temperature).

1. Draw a Bayesian network for this domain, given that the gauge is more likely to fail when the core temperature gets too high.



2. From your network topology, can you conclude $T \perp F_A \mid F_G$ and $T \perp F_A \mid F_G, A$?

All the path from T to F_A are (T, F_G, G, A, F_A) and (T, G, A, F_A) .

Knowing F_G , the triplet (G, A, F_A) is inactive (v -structure). Since it is present in both path, we can conclude that there are no active path, so $T \perp F_A \mid F_G$ is true.

For the second relation, the triplet (G, A, F_A) becomes active as we observe A . The only way we can ensure conditional independence is by ensuring that at least one remaining triplet is inactive in each path. For the first path, we have (T, F_G, G) which is inactive (cascade) which makes the path inactive. However, for the second path, (T, G, A) is active (cascade) which makes the path active. We have at least one active path from T to F_A which prevents us to conclude the conditional independence in this case.

3. Suppose there are just two possible actual and measured temperatures: low (l) and high (h). The probability that the gauge gives the correct temperature is x when it is working, but y when it is faulty. Give the conditional probability table associated with G .

The CPT for G is shown below. Students should pay careful attention to the semantics of F_G , which is true when the gauge is faulty, *i.e.* not working properly.

T	F_G	$P(G = l \mid T, F_G)$
l	0	x
l	1	y
h	0	$1 - x$
h	1	$1 - y$

4. Suppose the alarm is always triggered by high measured temperatures, unless it is faulty, in which case it never sounds. Give the conditional probability table associated with A .

G	F_A	$P(A = 1 G, F_A)$
l	0	0
l	1	0
h	0	1
h	1	0

5. Suppose the gauge is not faulty and the alarm is triggered. Calculate an expression for the probability that the temperature of the core is too high, in terms of the various conditional probabilities in the network.

The probability of interest here is $P(T = h|A = 1, F_G = 0)$. Because the alarm's behavior is deterministic, we can deduce that G indicates high (h). Hence, $P(T = h|A = 1, F_G = 0) = P(T = h|A = 1, G = h, F_G = 0)$. We also see that $A \perp T|G$. Therefore, we only need to calculate

$$\begin{aligned}
 P(T = h|G = h, F_G = 0) &= \frac{P(T = h, G = h, F_G = 0)}{P(G = h, F_G = 0)} \\
 &= \frac{P(G = h|T = h, F_G = 0)P(F_G = 0|T = h)P(T = h)}{\sum_t P(G = h|t, F_G = 0)P(F_G = 0|t)P(t)},
 \end{aligned}$$

which we cannot develop more as we don't know $P(T)$ and $P(F_G|T)$.

Exercise 2 Student–TA relationship

The teaching assistant of the course is worried about the time he spends to prepare the exercise sessions. He wants to investigate the influence of students can have on his schedule. To do so, he builds the following model



where S is a r.v. denoting the number of students leaving the room between a theoretical lecture and the practical session that follows, and H is a r.v. representing the number of hours spent to prepare the next session. He chooses the following relations for each variable

$$P(S = s) = \text{Poisson}(S = s; \lambda) = \frac{\lambda^s e^{-\lambda}}{s!} \quad (1)$$

$$H = \omega S + H_0 + \mathcal{N}(0, \sigma^2) \quad (2)$$

where H_0 corresponds to the number of hours the TA should spend per session regarding its contract and ω is the influence weight of students over the TA.

1. Identify the parameters of the Bayesian network and give the joint log-likelihood of one pair (s_i, h_i) given the model.

Given the topology of the network, we have to identify $P(S = s)$ and $P(H = h \mid S = s)$. The former is given by (1) whereas the latter can be deduced from (2). We have

$$P(S = s) = \text{Poisson}(S = s; \lambda)$$

and

$$P(H = h \mid S = s) = \mathcal{N}(\omega s + H_0, \sigma^2).$$

The parameters of those distributions are $\theta = [\lambda, \omega, \sigma]$.

The likelihood of a pair (s_i, h_i) can be derived from the network. We have

$$\begin{aligned} P((s_i, h_i) \mid \theta) &= P(s_i)P(h_i \mid s_i) \\ &= \frac{\lambda^{s_i} e^{-\lambda}}{s_i!} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(h_i - \omega s_i - H_0)^2}{2\sigma^2}} \end{aligned}$$

from which we derive the log-likelihood

$$LL((s_i, h_i); \theta) = s_i \log \lambda - \lambda - \log s_i! - \log \sqrt{2\pi}\sigma - \frac{(h_i - \omega s_i - H_0)^2}{2\sigma^2}. \quad (3)$$

2. The TA collects a set $\mathcal{D} = \{(s_i, h_i)\}_{i=1}^N$ of independent observations. Derive the maximum likelihood estimate of the parameters.

We first have to define the joint likelihood of those observations. Since they are assumed i.i.d., we have

$$P(\mathcal{D} \mid \theta) = \prod_{i=1}^N P((s_i, h_i) \mid \theta)$$

Hence, deriving the log-likelihood

$$\begin{aligned} LL(\mathcal{D}; \theta) &= \sum_i LL((s_i, h_i); \theta) \\ &= -N(\lambda + \log \sqrt{2\pi}\sigma) + \sum_i s_i \log \lambda - \log s_i! - \frac{(h_i - \omega s_i - H_0)^2}{2\sigma^2} \end{aligned}$$

The maximum likelihood parameters can be found by cancelling the derivative of the log-likelihood w.r.t. each parameter.

For λ we have

$$\begin{aligned}\frac{\partial LL}{\partial \lambda} &= -N + \sum_i \frac{s_i}{\lambda} \\ &= 0\end{aligned}$$

which gives $\lambda_{\text{MLE}} = \frac{1}{N} \sum_i s_i$.

For σ we have

$$\begin{aligned}\frac{\partial LL}{\partial \sigma} &= -\frac{1}{\sigma} \left(N - \sum_i \frac{(h_i - \omega s_i - H_0)^2}{\sigma^2} \right) \\ &\propto N - \frac{1}{\sigma^2} \sum_i (h_i - \omega s_i - H_0)^2 \\ &= 0\end{aligned}$$

which gives $\sigma^2_{\text{MLE}} = \frac{1}{N} \sum_i (h_i - \omega s_i - H_0)^2$.

For ω we have

$$\begin{aligned}\frac{\partial LL}{\partial \omega} &= \sum_i \frac{s_i (h_i - \omega s_i - H_0)}{\sigma^2} \\ &\propto \left(\sum_i s_i (h_i - H_0) \right) - \omega \sum_i s_i^2 \\ &= 0\end{aligned}$$

which gives $\omega_{\text{MLE}} = \frac{\sum_i s_i (h_i - H_0)}{\sum_i s_i^2}$.

3. The TA has just started to collect observations of students leaving the classroom. Being an adept of the Bayes' school, he knows that using MLE can lead to overfitting. Hence, he decides to incorporate uncertainty in his analysis of the parameter λ as he has no strong knowledge about it. Using a prior $\lambda \sim \text{Gamma}(\alpha, \beta)$, what would be the posterior distribution $P(\lambda \mid \mathcal{S} = \{s_i\}_{i=1}^N)$?

Hint: $\text{Gamma}(x \mid \alpha, \beta) = \frac{1}{Z} x^{\alpha-1} e^{-\beta x}$, with Z a normalizing constant.

Using the Bayes' theorem

$$\begin{aligned}P(\lambda \mid \mathcal{S}) &\propto P(\mathcal{S} \mid \lambda) P(\lambda) \\ &\propto \lambda^{\alpha-1} e^{-\beta \lambda} \prod_i \frac{\lambda^{s_i} e^{-\lambda}}{s_i!} \\ &= \lambda^{\alpha-1} e^{-\beta \lambda} \left(\frac{\lambda^{\sum_i s_i} e^{-N\lambda}}{\prod_i s_i!} \right) \\ &\propto \lambda^{(\alpha + \sum_i s_i - 1)} e^{-(\beta + N)\lambda}\end{aligned}$$

which corresponds to $\text{Gamma}(\alpha + \sum_i s_i, \beta + N)$, the posterior density. Note that \propto stands for the normalizing constant of the posterior density. Everything that does not depend explicitly on the parameter λ can be absorbed in the constant if it multiplies the whole.

4. Knowing that the mean, the variance and the mode of $\text{Gamma}(\alpha, \beta)$ are respectively $\frac{\alpha}{\beta}$, $\frac{\alpha}{\beta^2}$ and $\frac{\alpha - 1}{\beta}$, interpret how they evolve w.r.t. N between the prior and the posterior.

We have

$$\mathbb{E}(\lambda) : \frac{\alpha}{\beta} \Rightarrow \frac{\alpha + \sum_i s_i}{\beta + N},$$

$$\mathbb{V}(\lambda) : \frac{\alpha}{\beta^2} \Rightarrow \frac{\alpha + \sum_i s_i}{(\beta + N)^2}$$

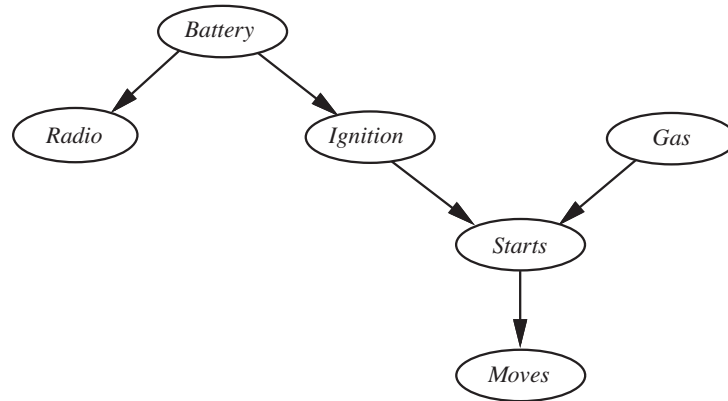
and

$$\arg \max_{\lambda} P(\lambda | \cdot) : \frac{\alpha - 1}{\beta} \Rightarrow \frac{\alpha + \sum_i s_i - 1}{\beta + N}.$$

When $N = 0$, the prior and the posterior are similar, which is expected. When $N \rightarrow \infty$, knowing that $s \in \mathbb{N}$ by definition (number of students leaving the room), we can assume that, for reasonable values of α, β , we have $\sum_i s_i \gg \alpha$ and $N \gg \beta$. This implies that the mean and the mode converge to the same value corresponding to λ_{MLE} . The variance decreases towards 0, which suggests that λ_{MLE} is the best and unique estimate of λ for an infinite amount of observations. This illustrates the trade-off between likelihood and prior in a Bayesian parameter learning setting. As more data are observed, the uncertainty about a parameter decreases smoothly from the prior to the likelihood.

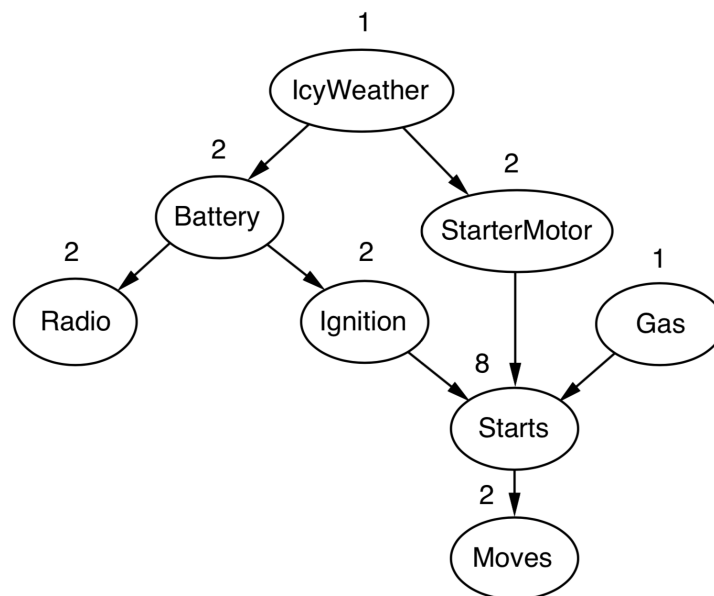
Exercise 3 Car diagnosis (AIMA, Ex 14.8)

Let be the following Bayesian network describing some features of a car's electrical system and engine. Each variable is Boolean, and the true value indicates that the corresponding aspect of the vehicle is in working order.



1. Extend the network with the Boolean variables *IcyWeather* and *StarterMotor*.

IcyWeather is not caused by any of the car-related variables, so needs no parents. It directly affects the battery and the starter motor. *StarterMotor* is an additional precondition for *Starts*.



2. According to your knowledge of cars, give reasonable conditional probability tables for all the nodes.

Reasonable probabilities may vary a lot depending on the kind of car and perhaps the personal experience of the assessor. The following values indicate the general order of magnitude and relative values that would be reasonable:

- A reasonable prior for *IcyWeather* might be 0.05 (depending on location and season).
- $P(\text{Battery}|\text{IcyWeather}) = 0.95, P(\text{Battery}|\neg\text{IcyWeather}) = 0.997$.
- $P(\text{StarterMotor}|\text{IcyWeather}) = 0.98, P(\text{StarterMotor}|\neg\text{IcyWeather}) = 0.999$.
- $P(\text{Radio}|\text{Battery}) = 0.9999, P(\text{Radio}|\neg\text{Battery}) = 0.05$.

- $P(\text{Ignition}|\text{Battery}) = 0.998, P(\text{Ignition}|\neg\text{Battery}) = 0.01.$
- $P(\text{Gas}) = 0.995.$
- $P(\text{Starts}|\text{Ignition}, \text{StarterMotor}, \text{Gas}) = 0.9999,$ other entries 0.0.
- $P(\text{Moves}|\text{Starts}) = 0.998.$

3. How many independent values are contained in the joint probability distribution for eight Boolean nodes, assuming that no conditional independence relations are known to hold among them?

With 8 Boolean variables, the joint has $2^8 - 1 = 255$ independent entries.

4. How many independent probability values do your network tables contain?

Given the topology with IcyWeather and StarterMotor, the total number of independent CPT entries is $1 + 2 + 2 + 2 + 2 + 1 + 8 + 2 = 20.$

Exercise 4 Gaussian parameter learning with known σ^2

You have in your possession a high-tech laser that emits light of intensity μ_l with great precision. You want to use it to test a new light-intensity measurement device. You assume that your measures y_i follow a Gaussian distribution $\mathcal{N}(\mu_y, \sigma^2)$ with $\mu_y \sim \mathcal{N}(\mu_l, \sigma_l^2)$ corresponding to the model of your laser given in the manual. You assume that σ^2 is fixed and that you can determine it.

You have access to \mathcal{Y}_c a calibration set of C observations and you want to make N new measures with your device. You want to study how you can reduce your uncertainty about the measures using \mathcal{Y}_c .

1. Compute the prior and posterior predictive distributions $P(y)$ and $P(y | \mathcal{Y}_c)$ if you have access to one observation, i.e. $\mathcal{Y}_c = \{y_o\}$.

The following exercise uses the Gaussian model identities (cfr. theoretical reminder of this session).

We first identify the components we have:

- $P(y | \mu_y) = \mathcal{N}(\mu_y, \sigma^2)$
- $P(\mu_y) = \mathcal{N}(\mu_l, \sigma_l^2)$

Prior predictive

$$P(y) = \int P(y | \mu_y)P(\mu_y)d\mu_y.$$

From Gaussian identities, you can identify different parameters as follows

$$\begin{aligned}\mathbf{x} &= \mu_y \\ \mathbf{m} &= \mu_l \\ \mathbf{P} &= \begin{pmatrix} \sigma_l^2 & 0 \\ 0 & \sigma_l^2 + \sigma^2 \end{pmatrix} \\ \mathbf{y} &= y \\ \mathbf{H} &= 1 \\ \mathbf{u} &= 0 \\ \mathbf{R} &= \sigma^2\end{aligned}$$

which corresponds to the joint

$$P(\mu_y, y) = \mathcal{N} \left(\begin{pmatrix} \mu_l \\ \mu_l \end{pmatrix}, \begin{pmatrix} \sigma_l^2 & \sigma_l^2 \\ \sigma_l^2 & \sigma_l^2 + \sigma^2 \end{pmatrix} \right)$$

from which we extract $P(y) = \mathcal{N}(\mu_l, \sigma^2 + \sigma_l^2)$.

Posterior predictive

$$P(y | \mathcal{Y}_c) = \int P(y | \mu_y)P(\mu_y | \mathcal{Y}_c)d\mu_y.$$

We first have to identify the posterior $P(\mu_y | \mathcal{Y}_c)$. Once again, from Gaussian identities,

knowing $P(\mathcal{Y}_c | \mu_y)$ and $P(\mu_y)$ you can identify different parameters as follows

$$\begin{aligned}\mathbf{x} &= \mu_y \\ \mathbf{m} &= \mu_l \\ \mathbf{P} &= \sigma_l^2 \\ \mathbf{y} &= y_o \\ \mathbf{H} &= 1 \\ \mathbf{u} &= 0 \\ \mathbf{R} &= \sigma^2\end{aligned}$$

which corresponds to the joint

$$P(\mu_y, y_o) = \mathcal{N}\left(\begin{pmatrix} \mu_l \\ \mu_l \end{pmatrix}, \begin{pmatrix} \sigma_l^2 & \sigma_l^2 \\ \sigma_l^2 & \sigma_l^2 + \sigma^2 \end{pmatrix}\right)$$

from which we extract

$$\begin{aligned}P(\mu_y | y_o) &= \mathcal{N}(\mu_{post}, \sigma_{post}^2) \\ &= \mathcal{N}\left(\mu_l + \frac{\sigma_l^2}{\sigma^2 + \sigma_l^2}(y_o - \mu_l), \sigma_l^2 - \frac{\sigma_l^4}{\sigma^2 + \sigma_l^2}\right).\end{aligned}$$

After simplification, we identify the posterior parameters as

$$\mu_{post} = \frac{\sigma^2}{\sigma^2 + \sigma_l^2}\mu_l + \frac{\sigma_l^2}{\sigma^2 + \sigma_l^2}y_o$$

and

$$\sigma_{post}^2 = \frac{\sigma^2}{\sigma_l^2 + \sigma^2}\sigma_l^2.$$

Following the same development as for the prior predictive, we have

$$P(y | \mathcal{Y}_c) = \mathcal{N}(\mu_{post}, \sigma^2 + \sigma_{post}^2).$$

2. Looking at you predictive uncertainty, what is the benefit of incorporating \mathcal{Y}_c ?

As a reminder, we found

$$\begin{aligned}P(y) &= \mathcal{N}(\mu_l, \sigma^2 + \sigma_l^2). \\ P(y | \mathcal{Y}_c) &= \mathcal{N}(\mu_{post}, \sigma^2 + \sigma_{post}^2).\end{aligned}$$

Comparing only the variances, we have

$$\begin{aligned}\sigma^2 + \sigma_l^2 &\stackrel{?}{\leq} \sigma^2 + \sigma_{post}^2 \\ \sigma_l^2 &\stackrel{?}{\leq} \sigma_{post}^2 \\ 1 &\stackrel{?}{\leq} \frac{\sigma^2}{\sigma_l^2 + \sigma^2} \\ \sigma_l^2 &\geq 0\end{aligned}$$

and we deduce that the posterior predictive has a smaller or equal variance to the prior one. We conclude that incorporating the calibration set can reduce uncertainty about the

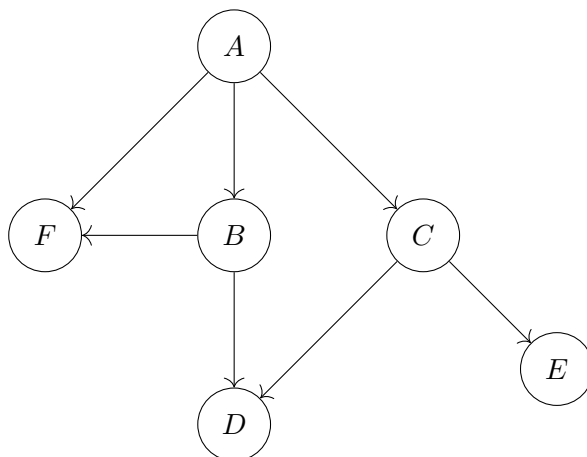
measures.

If we computed the variance for a calibration set of C observations, we would have found a posterior variance

$$\sigma_{post}^2 = \frac{\sigma^2}{C\sigma_l^2 + \sigma^2}\sigma_l^2.$$

Since σ^2 is fixed, we see that including more observations reduces even more the posterior predictive variance. In the limit case $C \rightarrow \infty$, the posterior predictive variance converges to σ^2 which corresponds to the uncertainty of the instrument (you can't reach a lower uncertainty, it is called the aleatoric uncertainty).

Exercise 5 Independence



Considering the hereabove Bayesian network, which of the following statements are enforced by the network structure?

We apply the d-separation algorithm. To show that two variables X and Y could be dependent, it is sufficient to find a single *active* (undirected) path from X to Y . A path is active if all of its consecutive triplets are active.

1. $P(B, C) = P(B)P(C)$

This is true iff (if and only if) $B \perp C$ (B is independent from C). The three paths that link B and C are (B, A, C) , (B, D, C) and (B, F, A, C) . (B, A, C) is active because A is unknown (\wedge -structure). We cannot guarantee that $B \perp C$.

2. $P(B, C|A) = P(B|A)P(C|A)$

This is true iff $B \perp C|A$. This time, (B, A, C) is not active because A is known. (B, D, C) is not active either because D is unknown (\vee -structure). In (B, F, A, C) , the triplet (B, F, A) is active, but (F, A, C) is not. We can guarantee that $B \perp C|A$.

3. $P(F, E) = P(F)P(E)$

This is true iff $F \perp E$. The four paths that link F and E are (F, A, C, E) , (F, A, B, D, C, E) , (F, B, A, C, E) and (F, B, D, C, E) . In (F, A, C, E) , the triplets (F, A, C) (\wedge -structure) and (A, C, E) are both active. We cannot guarantee that $F \perp E$.

4. $P(F|A, D, E) = P(F|A, D)$

This is true iff $F \perp E|A, D$. This time, the triplets (F, A, C) , (F, A, B) and (B, A, C) are inactive (\wedge -structure, but A is known). In (F, B, D, C, E) , (F, B, D) is active. Indeed, knowing an extremity of a triplet does not impact whether it is active or not. Only the center variable is important. Then we have (B, D, C) which is active (\vee -structure, but D is known) and (D, C, E) also. Hence, (F, B, D, C, E) is active and we cannot guarantee that $F \perp E|A, D$.

5. $P(B, E) = \sum_{a,c,d,f} P(a)P(B|a)P(c|a)P(d|B, c)P(E|c)P(f|a, B)$

We notice that

$$P(A, B, C, D, E, F) = P(A)P(B|A)P(C|A)P(D|B, C)P(E|C)P(F|A, B)$$

corresponds exactly to

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | \text{parents}(X_i)),$$

for the considered network, which we know is guaranteed. We also know that

$$P(B, E) = \sum_{a,c,d,f} P(a, B, c, d, E, f)$$

Hence, the initial statement is always true.

For the same network, use inference by variable elimination to compute $P(E|A = 1, B = 1)$.

We have

$$\begin{aligned} P(E|A = 1, B = 1) &= \alpha \sum_{c,d} P(A = 1, B = 1, c, d, E) \\ &= \alpha \sum_{c,d} P(A = 1)P(B = 1|A = 1)P(c|A = 1)P(d|B = 1, c)P(E|c) \\ &= \alpha P(A = 1)P(B = 1|A = 1) \sum_c P(c|A = 1)P(E|c) \sum_d P(d|B = 1, c). \end{aligned}$$

We define the initial factors as

$$\begin{aligned} f_1 &= P(A = 1) \\ f_2 &= P(B = 1|A = 1) \\ f_3(C) &= P(C|A = 1) \\ f_4(E, C) &= P(E|C) \\ f_5(C, D) &= P(D|B = 1, C) \end{aligned}$$

and the composite factors as

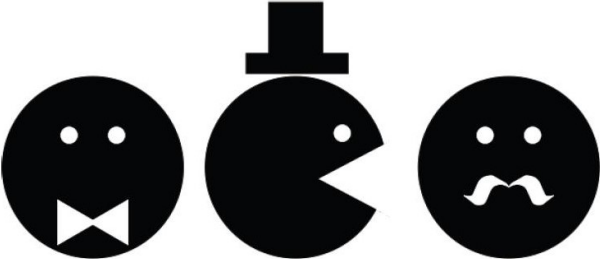
$$\begin{aligned} f_6(C) &= \sum_d f_5(C, d) \\ f_7(C, E) &= f_3(C) \times f_4(E, C) \times f_6(C) \\ f_8(E) &= \sum_c f_7(c, E) \\ f_9(E) &= f_1 \times f_2 \times f_7(E). \end{aligned}$$

Finally,

$$P(E|A = 1, B = 1) = \alpha f_9(E) = \frac{f_9(E)}{\sum_e f_9(e)}.$$

Exercise 6 Pacbaby (UC Berkeley CS188, Spring 2014)

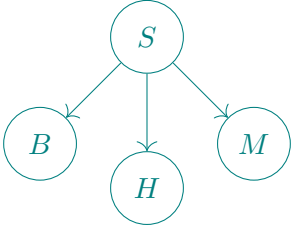
Pacman and Pacwoman have been searching for each other in the maze. Pacwoman has been pregnant with a baby, and just this morning she has given birth to Pacbaby¹. Because Pacbaby was born before Pacman and Pacwoman were reunited in the maze, he has never met his father. Naturally, Pacwoman wants to teach Pacbaby to recognize his father, using a set of pictures of Pacman. She also has several pictures of ghosts to use as negative examples.



Because the pictures are black and white, and were taken from various angles, Pacwoman has decided to teach Pacbaby to identify Pacman based on salient features: the presence of a bowtie B , hat H or mustache M . The following table summarizes the content of the pictures. Each feature takes realization in $\{0, 1\}$, where 0 and 1 mean the feature is respectively absent and present. The subject of the picture is described by a random variable $S \in \{0, 1\}$, where 0 is a ghost and 1 is Pacman.

B	H	M	S
0	0	0	1
1	0	0	0
1	0	1	1
1	1	0	0
0	1	0	0
1	1	1	1

1. Suppose Pacbaby has a Naive Bayes based brain. Draw the Bayesian network that would represent the dependencies between S , B , H and M for Pacbaby.



2. Write the Bayesian classification rule for this problem, *i.e.* the formula that given a data point (b, h, m) returns the most likely subject. Write the formula in terms of conditional and prior probabilities. What does the formula become under the assumptions of Pacbaby?

Given (b, h, m) , the most likely subject is given by the *maximum a posteriori* (MAP)

¹Congratulations!

estimation

$$\begin{aligned} s_{\text{MAP}} &= \arg \max_s P(s|b, h, m) \\ &= \arg \max_s P(b, h, m|s)P(s). \end{aligned}$$

Under the naive Bayes assumptions of Pacbaby, B , H and M become independent conditionally to S , *i.e.* $P(B, H, M|S) = P(B|S)P(H|S)P(M|S)$. Then, the formula becomes

$$s_{\text{MAP}} = \arg \max_s P(b|s)P(h|s)P(m|s)P(s).$$

3. What are the parameters of this model? Give estimates of these parameters according to the pictures provided by Pacwoman.

The parameters of the model are the elements of the prior vector $P(S)$ and the (conditional) probability matrices $P(B|S)$, $P(H|S)$ and $P(M|S)$. An (unbiased) estimation of these elements can be computed as the frequency of their respective events within the learning set (of pictures).

S	$P(S)$	$P(B = 1 S)$	$P(H = 1 S)$	$P(M = 1 S)$
0	$\frac{3}{6}$	$\frac{2}{3}$	$\frac{2}{3}$	$\frac{0}{3}$
1	$\frac{3}{6}$	$\frac{2}{3}$	$\frac{1}{3}$	$\frac{2}{3}$

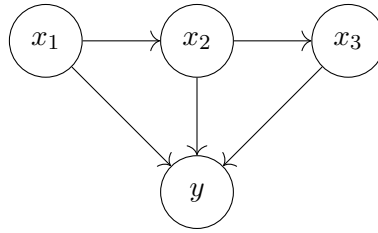
4. Pacman eventually shows up wearing a bowtie, but no hat or mustache. Will Pacbaby recognize his father?

Pacbaby will recognize his father if $s_{\text{MAP}} = 1$ for $(b, h, m) = (1, 0, 0)$. Using the parameters estimated previously, we have

$$\begin{aligned} P(b|0)P(h|0)P(m|0)P(0) &= \frac{2}{3} \times \left(1 - \frac{2}{3}\right) \times \left(1 - \frac{0}{3}\right) \times \frac{3}{6} \approx 0.111 \\ P(b|1)P(h|1)P(m|1)P(1) &= \frac{2}{3} \times \left(1 - \frac{1}{3}\right) \times \left(1 - \frac{2}{3}\right) \times \frac{3}{6} \approx 0.074. \end{aligned}$$

Therefore, $s_{\text{MAP}} = 0$, meaning that Pacbaby will *not* recognize his father.

Exercise 7 Predict your grade



The hereabove Bayesian network represents how the final grade of a class is computed. In this model, x_1 , x_2 and x_3 respectively denote the grades obtained by a student at the homework, project and exam. The teaching assistant that grades the homework also grades the project and the exam, which introduces a slight bias in the corrections. In particular, $x_2 \sim \mathcal{N}(a_1x_1 + \mu_2, \sigma_2^2)$ and $x_3 \sim \mathcal{N}(a_2x_2 + \mu_3, \sigma_3^2)$. Finally, $y \sim \mathcal{N}(a_3x_1 + a_4x_2 + a_5x_3 + \mu_y, \sigma_y^2)$ stands for the final grade, which is a linear combination of the grades obtained by the student during the semester plus some Gaussian noise due to rounding errors. Answer the following questions about this model.

1. Assuming the parameters of the model are known, what is the expected value of y given x_1 and x_2 .

Our task is to find the expectation

$$\mathbb{E}_{p(y|x_1, x_2)}[y] = \int y p(y|x_1, x_2) dy.$$

We know that

$$p(y|x_1, x_2) = \int p(y|x_1, x_2, x_3) p(x_3|x_1, x_2) dx_3,$$

where $p(y|x_1, x_2, x_3)$ and $p(x_3|x_1, x_2)$ are linear Gaussian distributions given in the statement. Therefore, we have

$$p(y|x_1, x_2) = \mathcal{N}\left(a_3x_1 + a_4x_2 + a_5(a_2x_2 + \mu_3) + \mu_y, (a_5\sigma_3)^2 + \sigma_y^2\right)$$

and, by definition of a Gaussian distribution,

$$\mathbb{E}_{p(y|x_1, x_2)}[y] = a_3x_1 + (a_4 + a_5a_2)x_2 + a_5\mu_3 + \mu_y.$$

2. Suppose now that the model's parameters are unknown. Given a learning set $d = \{(x_{i,1}, x_{i,2}, y_i)\}$ of N independent and identically distributed points, determine the model that best describes d .

We know that the distribution of y given x_1 and x_2 takes the form $\mathcal{N}(w_1x_1 + w_2x_2 + b, \sigma^2)$. Then, our task is to find the parameters $h = (w_1, w_2, b, \sigma)$ that maximize the likelihood of d , i.e. the *maximum likelihood estimation* (MLE)

$$\begin{aligned} h_{\text{MLE}} &= \arg \max_w p(d|h) \\ &= \arg \max_h \prod_i p(x_i, y_i|h) \\ &= \arg \max_h \log \prod_i p(x_i, y_i|h) \\ &= \arg \max_h \sum_i \log p(x_i, y_i|h) \end{aligned}$$

$$\begin{aligned}
&= \arg \max_h \sum_i \log p(y_i|h, x_i) + \log p(x_i) \\
&= \arg \max_h \sum_i \log p(y_i|h, x_i) \\
&= \arg \max_h \sum_i \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(w^T x_i - y_i)^2}{2\sigma^2}\right) \right] \\
&= \arg \max_h \sum_i -\frac{1}{2} \log(2\pi\sigma^2) - \frac{(w^T x_i - y_i)^2}{2\sigma^2} \\
&= \arg \min_h \log \sigma^2 + \frac{1}{\sigma^2} \frac{1}{N} \sum_i (w^T x_i - y_i)^2,
\end{aligned}$$

where $x_i = (x_{i,1} \ x_{i,2} \ 1)^T$ and $w = (w_1 \ w_2 \ b)^T$. In the last expression, we observe that the summation term is independent from σ . Therefore,

$$w_{\text{MLE}} = \arg \min_w \sum_i (w^T x_i - y_i)^2,$$

which exactly corresponds to a *linear regression* problem. Then, we find w_{MLE} by canceling the gradient with respect to w , *i.e.*

$$\begin{aligned}
0 &= \nabla_w \sum_i (w^T x_i - y_i)^2 \\
&= \nabla_w \sum_i (w^T x_i - y_i)(w^T x_i - y_i) \\
&= \nabla_w \sum_i (w^T x_i)^2 + y_i^2 - 2w^T x_i y_i \\
&= \nabla_w (w^T X^T X w + Y^T Y - 2w^T X^T Y) \\
&= 2X^T X w - 2X^T Y
\end{aligned}$$

where $X = (x_i^T) \in \mathbb{R}^{N \times 3}$ and $Y = (y_i) \in \mathbb{R}^N$. Finally, we have

$$\begin{aligned}
0 &= X^T X w_{\text{MLE}} - X^T Y \\
\Leftrightarrow w_{\text{MLE}} &= (X^T X)^{-1} X^T Y.
\end{aligned}$$

Afterwards, we find σ_{MLE} such that

$$\begin{aligned}
\sigma_{\text{MLE}} &= \arg \min_{\sigma} \log \sigma^2 + \frac{\text{MSE}}{\sigma^2} \\
&= \sqrt{\text{MSE}},
\end{aligned}$$

where MSE denotes the *mean squared error*

$$\frac{1}{N} \sum_i (w_{\text{MLE}}^T x_i - y_i)^2.$$

Exercise 8 Heteroscedastic linear regression

What becomes the expression of the weight vector w in the solution of question 2.2 if the noise is different for each sample? In particular, $y_i \sim N(w^T x, \sigma_i^2)$ and we know the values σ_i .

Quiz

The Normal distribution $\mathcal{N}(\mu, \sigma)$ is described by the density function

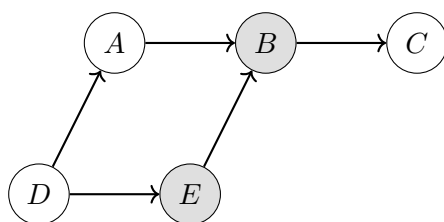
$p(x) = \frac{1}{\sigma} \exp(-z + \exp(-z))$, with $z = \frac{x-\mu}{\sigma}$.

$p(x) = \frac{1}{2\sigma} \exp\left(-\frac{|x-\mu|}{\sigma}\right)$.

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

$p(x) = 1 / \pi\sigma \left[1 + \left(\frac{x-\mu}{\sigma}\right)^2\right]$.

Consider the Bayesian network shown below. We want to infer $\mathbf{P}(A|b, e)$ where A is the query variable, B and E are evidence variables, and C and D are hidden variables. Which of the following statements is true?



$\mathbf{P}(A|b, e) \propto \sum_c \sum_d P(c|b)\mathbf{P}(b|A)P(b|e)\mathbf{P}(A|d)P(e|d)P(d)$

$\mathbf{P}(A|b, e) \propto \sum_c \sum_d P(c|b)\mathbf{P}(b|A, e)\mathbf{P}(A|d)P(e|d)$

$\mathbf{P}(A|b, e) \propto \sum_c \sum_d P(c|b)\mathbf{P}(b|A, e)\mathbf{P}(A|d)P(e|d)P(d)$

$\mathbf{P}(A|b, e) \propto \sum_b \sum_e P(c|b)\mathbf{P}(A|b, e)\mathbf{P}(A|d)P(e|d)P(d)$

In a Bayesian network

Independence between two variables is guaranteed if all path connecting them are inactive.

We can guarantee dependence between two variables using d-separation.

A path between two nodes is active if it contains at least one active triplet.

A cascade triplet is active if the center node is observed.

In a one node Bayesian network with a binary variable following a Bernoulli distribution. If we observe T positive realizations and F negative ones,

the maximum likelihood estimate of the positive probability is $\frac{T}{F}$.

the maximum likelihood estimate of the positive probability is $\frac{T+F}{T}$.

the maximum likelihood estimate of the positive probability is $\frac{T}{T+F}$.

None of the above.