

1

$$(a) P(T, X | M, B) \stackrel{iid}{=} \prod_{n=1}^N P(\vec{t}_n, \vec{x}_n | M, B) =$$

$$\stackrel{\text{Prod}}{=} \prod_{n=1}^N p(\vec{x}_n | \vec{t}_n, M, B) \cdot P(\vec{t}_n | M, B) =$$

Rule t_{nk}

$$= \prod_{n=1}^N \prod_{k=1}^K \left[p(\vec{x}_n | t_{nk}=1, \vec{\mu}_k, \vec{\beta}_k) \cdot P(t_{nk}=1 | \vec{\mu}_k, \vec{\beta}_k) \right]$$

We set $t_{nk} = 1$ because we know that in the case where $t_{nk} = 0$ the exponent will be 0, which will cause the whole relation to be 1 (which doesn't affect the product). Also, we can say that $p(t_{nk}=1 | \vec{\mu}_k, \vec{\beta}_k) = p(C_k | \vec{\mu}_k, \vec{\beta}_k) = p(C_k) = \pi_{lk}$. Therefore, we continue with the Naive Bayes assumption:

$$\prod_{n=1}^N \prod_{k=1}^K \left[\pi_{lk} \cdot \prod_{d=1}^D p(x_{nd} | C_k, \vec{\mu}_{dk}, \vec{\beta}_{dk}) \right]^{t_{nk}}$$

C_k

(b) For a start, we need to compute $P(X_{n1}, X_{n2}, \dots, X_{nD} | \vec{\mu}_1, \vec{\beta}_1)$, which is the joint distribution of all the features. With Naive Bayes, we assume that they are independent, which means that we don't take into consideration potential correlation between them, we assume that there is none. In this case, we have D · K parameters, as we have D variance vectors for each one of the K classes, and D mean vectors for each K class. On the other hand, if we don't assume independence we should take into account of all the potential correlations between the features (in fact, this is the most significant benefit of Naive Bayes, as this is $O(2^D)$ and by assuming independence this is reduced to $O(D)$). In this case, we get a D · D covariance

matrix for each of the K classes which is $D \cdot D \cdot K$ parameters, as well as $D \cdot K$ parameters for the K mean vectors of size D . Finally, since the covariance matrix is symmetric we need to subtract $(D+1)/2$ from the original D parameters. The final number of parameters thus is $D \cdot K \cdot (D - \frac{D+1}{2}) + D \cdot K$ which is equal to $DK \cdot \binom{D+3}{2}$. Furthermore, an example where this assumption would lead to incorrect results, is an example where the features are correlated. Assume that we are modelling car sales in 2024 and some of the features are ^{fuel}~~gas~~ type (electric, gas, petrol) and date that the car was released. The past years the released cars are mostly electric, in comparison to 10 years ago, therefore the features are not independent, and the assumption does not hold.

$$\begin{aligned}
 (c) \log \left(P(T, X | M, B) \right) &= \log \left(\prod_{n=1}^N \prod_{k=1}^K \left[\pi_k \cdot \prod_{d=1}^D p(x_{nd} | c_k, f_{dk}, \hat{p}_{dk}) \right]^{t_{nk}} \right) = \\
 &= \sum_{n=1}^N \sum_{k=1}^K t_{nk} \cdot \left(\log \pi_k + \sum_{d=1}^D \log p(x_{nd} | c_k, f_{dk}, \hat{p}_{dk}) \right) = \\
 &= \sum_{n=1}^N \sum_{k=1}^K t_{nk} \cdot \left(\log \pi_k + \sum_{d=1}^D \log \left(\frac{\hat{p}_{dk}^{1/2}}{2\pi^{1/2}} \cdot \exp \left\{ -\frac{\hat{p}_{dk}}{2} (x_{nd} - f_{dk})^2 \right\} \right) \right) = \\
 &= \sum_{n=1}^N \sum_{k=1}^K t_{nk} \cdot \left(\log \pi_k + \sum_{d=1}^D \left(\log(\hat{p}_{dk}^{1/2}) - \log(2\pi^{1/2}) - \frac{\hat{p}_{dk}}{2} (x_{nd} - f_{dk})^2 \right) \right) =
 \end{aligned}$$

This is as far as we can go so we move on with (d) and the derivation of the log likelihood relation that we ended up with.

(d) We continue and use the relation we derived in (c). We will compute its partial derivative in terms of an μ_{ij} th element and subsequently we will set it to 0 and solve for μ_{ij} .

$$\frac{\partial}{\partial \mu_{ij}} \left(\sum_{n=1}^N \sum_{k=1}^K t_{nk} \left(\log p_{ik} + \sum_{d=1}^D -\bar{p}_{dk} \cdot (x_{nd} - \mu_{dk})^2 \right) \right) = \frac{\log(\bar{p}_{ik}) - \log(x_{nd}^{-1})}{\text{discard the constants}}$$

$$= \sum_{n=1}^N \sum_{k=1}^K t_{nk} \cdot \sum_{d=1}^D -\bar{p}_{dk} \cdot \left(\frac{\partial}{\partial \mu_{ij}} (\bar{x}_{nd}) - \frac{\partial}{\partial \mu_{ij}} (2x_{nd} \mu_{dk}) + \frac{\partial}{\partial \mu_{ij}} (\mu_{dk}^2) \right) =$$

$$= \sum_{n=1}^N \sum_{k=1}^K t_{nk} \sum_{d=1}^D -\bar{p}_{dk} \cdot (-2x_{nd} \delta_{id} \delta_{jk} + 2\mu_{dk} \delta_{id} \delta_{jk}) =$$

Remove the parameter deltas as well as the corresponding summations and set $i=d=i$ and $k=j$.

$$\sum_{n=1}^N t_{nj} \bar{p}_{ij} \cdot (2x_{ni} - 2\mu_{ij}) \quad (1)$$

Set (1) equal to 0 and solve for μ_{ij} :

$$\sum_{n=1}^N t_{nj} \bar{p}_{ij} (2x_{ni} - 2\mu_{ij}) = 0 \iff$$

$$\stackrel{(M_2)}{\iff} \sum_{n=1}^N t_{nj} \bar{p}_{ij} \cdot x_{ni} - t_{nj} \bar{p}_{ij} \mu_{ij} = 0 \iff$$

$$\stackrel{p_{ij} \neq 0}{\iff} \sum_{n=1}^N t_{nj} x_{ni} = \sum_{n=1}^N t_{nj} \mu_{ij} \iff$$

$$\iff \boxed{\mu_{ij} = \frac{\sum_{n=1}^N t_{nj} x_{ni}}{\sum_{n=1}^N t_{nj}}}$$

④ We have to assume that p_{ij} is not 0, because if it was 0 then we wouldn't have a valid probability distribution. That is because if it was 0 it would imply that we have infinite variance and in any meaningful probabilistic model we expect to have at least some certainty about our data.

So the MLE of μ_{dk} is:

$$\mu_{dk} = \frac{\sum_{n=1}^N t_{nk} \cdot x_{nd}}{\sum_{n=1}^N t_{nk}}$$

Which describes the following:
 t_{nk} : the number of regions that were classified as k (k th climate)
 $x_{nd} \cdot t_{nk}$: the number of measurements that occurred in the k th class

So we can also write the result as: $\mu_{dk} = \frac{N_{dk}}{N_K}$, which is the frequency of the regions that have the d th temperature measurement for the k th climate class.

$$\begin{aligned}
 (e) p(C_1 | x) &\stackrel{(1)}{=} P(x | C_1) \cdot p(C_1) \\
 &\stackrel{\text{P}}{=} \frac{\prod_{d=1}^D p(x_d | \mu_{d1}, \beta_{d1}) \cdot \pi_1}{\sum_{k=1}^K p(x | C_k) \cdot p(C_k)} = \\
 &= \frac{\prod_{d=1}^D p(x_d | \mu_{d1}, \beta_{d1}) \cdot \pi_1}{\sum_{k=1}^K \pi_k \cdot p(x | C_k)} = \\
 &\stackrel{(2)}{=} \frac{\prod_{d=1}^D p(x_d | \mu_{d1}, \beta_{d1}) \cdot \pi_1}{\sum_{k=1}^K \prod_{d=1}^D p(x_d | \mu_{dk}, \beta_{dk}) \cdot \pi_k} =
 \end{aligned}$$

⊕ In the second step we just marginalize over k using the sum rule in order to index M, B , as the exercise demands a result in terms of $p(x | \mu, B)$

⊗ In the first step as well as in the fourth step we replace C_k with μ_k, β_k as it is parameterized by these values so we can proceed with them.

(f) A generative model aims to model the joint distribution $p(x, c)$ which is the probability of both the features x and the class c . Unlike discriminative models, such as logistic regression, which directly model the conditional probability $p(c|x)$, generative models do this by factoring the joint distribution $p(x|c) \cdot p(c)$. Naive Bayes assumes that all features x_d are independent given the class c , which allows for the simplification of $p(x|c)$ as a product of individual feature likelihoods. This approach, where we first model the joint distribution and then (via Bayes' rule) we obtain $p(c|x)$ is what characterizes Naive Bayes as a generative model.

(g) We derived the following relation in (e) :

$$\frac{\prod_{d=1}^D P(x_d | \mu_{d1}, \beta_{d1}) \cdot \pi_1}{\sum_{k=1}^K \prod_{d=1}^D P(x_d | \mu_{dk}, \beta_{dk}) \cdot \pi_k} = \frac{\prod_{d=1}^D \left[\frac{\beta_{d1}^{1/2}}{2\pi^{1/2}} \cdot \exp \left\{ -\frac{\beta_{d1}}{2} (x_d - \mu_{d1})^2 \right\} \right] \pi_1}{\sum_{k=1}^K \prod_{d=1}^D \left[\frac{\beta_{dk}^{1/2}}{2\pi^{1/2}} \cdot \exp \left\{ -\frac{\beta_{dk}}{2} (x_d - \mu_{dk})^2 \right\} \right] \pi_k}$$

(h) The condition is :

$$P(C_1 | x_u) > P(C_k | x_u), \text{ where } k \in \mathbb{Z} \text{ and } k \neq 1, K$$

(K : total classes)

The above relation shows that the input vector x_u will be classified as " C_1 ", due to the probability of C_1 being higher than all of the other corresponding class probabilities.

(i) We have: $p(C_1 | X_d) > p(C_k | X_d)$ $\xrightarrow{\text{Bayes}}$
 $\xleftarrow{\text{Rate}}$

$$p(C_1) \cdot p(X_d | C_1) > p(C_k) \cdot p(X_d | C_k) \quad (1)$$

We replace \log to (1) accordingly:

$$\begin{aligned} \prod_{d=1}^D \frac{\vec{P}_{d1}}{2\pi}^{1/2} \exp \left\{ -\frac{\vec{P}_{d1}}{2} (\vec{x}_d - \vec{\mu}_{d1})^2 \right\} &> \prod_{d=1}^D \frac{\vec{P}_{dk}}{2\pi}^{1/2} \exp \left\{ -\frac{\vec{P}_{dk}}{2} (\vec{x}_d - \vec{\mu}_{dk})^2 \right\} = \\ = \prod_{d=1}^D \vec{P}_{d1}^{1/2} \exp \left\{ -\frac{\vec{P}_{d1}}{2} (\vec{x}_d - \vec{\mu}_{d1})^2 \right\} &> \prod_{d=1}^D \vec{P}_{dk}^{1/2} \exp \left\{ -\frac{\vec{P}_{dk}}{2} (\vec{x}_d - \vec{\mu}_{dk})^2 \right\} = \end{aligned}$$

We log both sides:

$$\begin{aligned} = \log \prod_{d=1}^D \vec{P}_{d1}^{1/2} - \frac{\vec{P}_{d1}}{2} (\vec{x}_d - \vec{\mu}_{d1})^2 &> \log \prod_{d=1}^D \vec{P}_{dk}^{1/2} - \frac{\vec{P}_{dk}}{2} (\vec{x}_d - \vec{\mu}_{dk})^2 = \\ = \sum_{d=1}^D -\vec{P}_{d1} (\vec{x}_d - \vec{\mu}_{d1})^2 + \vec{P}_{dk} (\vec{x}_d - \vec{\mu}_{dk})^2 &> 2 \log \left(\frac{\vec{P}_{dk}}{\vec{P}_{d1}} \right) + \sum_{d=1}^D \log \frac{\vec{P}_{dk}}{\vec{P}_{d1}} \end{aligned}$$

We set $2 \log \left(\frac{\vec{P}_{dk}}{\vec{P}_{d1}} \right) + \sum_{d=1}^D \log \frac{\vec{P}_{dk}}{\vec{P}_{d1}} = C$ as these are constants and we get:

$$\sum_{d=1}^D \left(-\vec{P}_{d1} (\vec{x}_d - \vec{\mu}_{d1})^2 + \vec{P}_{dk} (\vec{x}_d - \vec{\mu}_{dk})^2 \right) > C$$

We collapse the summation:

$$\vec{P}_1 (\vec{x} - \vec{\mu}_1) + \vec{P}_k (\vec{x} - \vec{\mu}_k) > C$$

We set $\vec{B}_1 \cdot \vec{I} = \vec{B}_1$, $\vec{P}_k \cdot \vec{I} = \vec{B}_k$

$$\begin{aligned} -\vec{B}_1 (\vec{x} - \vec{\mu}_1) + \vec{B}_k (\vec{x} - \vec{\mu}_k)^2 &> C \Leftrightarrow \\ \cancel{-\vec{B}_1 \vec{x} + \vec{B}_1 \vec{\mu}_1 + \vec{B}_k \vec{x} - \vec{B}_k \vec{\mu}_k} - (\vec{x} - \vec{\mu}_1)^T \vec{B}_1 (\vec{x} - \vec{\mu}_1) + (\vec{x} - \vec{\mu}_k)^T \vec{B}_k (\vec{x} - \vec{\mu}_k) &> C \Leftrightarrow \\ -\vec{x}^T \vec{B}_1 \vec{x} + \vec{x}^T \vec{B}_1 \vec{\mu}_1 + \vec{\mu}_1^T \vec{B}_1 \vec{x} - \vec{\mu}_1^T \vec{B}_1 \vec{\mu}_1 + \vec{x}^T \vec{B}_k \vec{x} - \vec{x}^T \vec{B}_k \vec{\mu}_k + \vec{\mu}_k^T \vec{B}_k \vec{x} - \vec{\mu}_k^T \vec{B}_k \vec{\mu}_k &> C \\ \Leftrightarrow \vec{x}^T \vec{B}_k \vec{x} - \vec{x}^T \vec{B}_1 \vec{x} - 2 \vec{x}^T \vec{B}_k \vec{\mu}_k + 2 \vec{x}^T \vec{B}_1 \vec{\mu}_1 + \vec{\mu}_k^T \vec{B}_k \vec{\mu}_k - \vec{\mu}_1^T \vec{B}_1 \vec{\mu}_1 &> C' \end{aligned}$$

where $C' = C - \vec{\mu}_k^T \vec{B}_k \vec{\mu}_k + \vec{\mu}_1^T \vec{B}_1 \vec{\mu}_1$

So we have: $\vec{x}^T (2 \vec{B}_1 \vec{\mu}_1 - 2 \vec{B}_k \vec{\mu}_k) + \vec{x}^T (\vec{B}_k - \vec{B}_1) \vec{x} > C'$

which is saying ~~not~~ the same as (1) but written as an inequality in terms of the posterior of \log .

(k) Since $\beta_{dh} = \beta_{dk}$ (i) becomes:

$$x^T (2B_{f1} - 2B_{fk}) > c' \quad (1)$$

In order to prove this, we want to bring somehow the relation $x_0(1-\lambda) + \lambda x_1$ from (i) to (1). One simple way is to just set $x = x_0(1-\lambda) + \lambda x_1$. Then (1) transforms to:

$$(x_0(1-\lambda) + \lambda x_1)^T (2B_{f1} - 2B_{fk}) > c'$$

We do the product and we get:

$$\underbrace{x_0^T(1-\lambda)}_{\lambda} \underbrace{(2B_{f1} - 2B_{fk})}_{+ \lambda x_1^T(2B_{f1} - 2B_{fk})} > c'$$

This is a convex combination of the decision inequalities at x_0 and x_1 , which accompanied by the assumption that the x_0, x_1 we selected are in the decision region C_1 , prove that any point x_λ on the straight line between x_0 and x_1 also satisfies the decision boundary condition for C_1 . Thus, the region is convex.

(j) If the quadratic term of (i) is convex, then the decision regions for this problem will be convex as well. The term is:

$$x^T (B_k - B_1) x \quad (1)$$

This term will be convex iff the matrix $B_k - B_1$ is positive definite, or $B_k = B_1$. The exercise hasn't specified the latter so we have to assume that it is $B_k \neq B_1$. Thus we need to prove that $B_k - B_1$ is positive definite in order to say that it is convex. Given the nature of ~~the problem~~ ^{matrix theory} it is very hard to be certain about the positive definiteness of $B_k - B_1$ because we are not sure that all eigenvalues will be larger than 0.

To conclude, it could be convex if, we knew $B_k = B_1$ (because then we have a linear problem as the quadratic term vanishes) or in the case where $B_k - B_1$ is positive definite. Since we don't know either, we can't answer with certainty. By intuition, if $B_k \neq B_1$ it will probably be non-convex, as some eigenvalues won't be positive.

(3) (a) We start by writing down the likelihood:

$$\begin{aligned} P(T | \Phi, w_1, \dots, w_k) &= \prod_{n=1}^N P(\vec{t}_n | \vec{\Phi}_n, w_1, \dots, w_k) \stackrel{\text{take is the}}{=} \text{same as } c_k \\ &= \prod_{n=1}^N \prod_{k=1}^K P(c_k | \vec{\Phi}_n)^{t_{nk}} \end{aligned}$$

Then, we derive the log-likelihood:

$$\begin{aligned} \log(P(T | \Phi, w_1, \dots, w_k)) &= \log\left(\prod_{n=1}^N \prod_{k=1}^K P(c_k | \vec{\Phi}_n)^{t_{nk}}\right) = \\ &= \sum_{n=1}^N \sum_{k=1}^K t_{nk} \cdot \log(P(c_k | \vec{\Phi}_n)) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \cdot \log(y_k(\vec{\Phi}_n)) \end{aligned}$$

$$\begin{aligned} (b) P(w_1, \dots, w_k | \alpha) &= \prod_{k=1}^K N(w_k | 0, \alpha^{-1} I) = \\ &= \prod_{k=1}^K 2\pi^{-\frac{k}{2}} \cdot |\alpha^{-1} \cdot I|^{-\frac{1}{2}} \cdot \exp\left\{-\frac{1}{2} \cdot w_k^\top (\alpha^{-1} \cdot I) \cdot w_k\right\} = \\ &\stackrel{\alpha^{-1} = I}{=} \prod_{k=1}^K 2\pi^{-\frac{k}{2}} \cdot |\alpha^{-1} \cdot I|^{-\frac{1}{2}} \cdot \exp\left\{-\frac{\alpha}{2} \cdot w_k^\top w_k\right\} \end{aligned}$$

Then, we compute the logarithm:

$$\begin{aligned} \log(P(w_1, \dots, w_k | \alpha)) &= \log\left(\prod_{k=1}^K 2\pi^{-\frac{k}{2}} \cdot |\alpha^{-1} \cdot I|^{-\frac{1}{2}} \cdot \exp\left\{-\frac{\alpha}{2} \cdot w_k^\top w_k\right\}\right) = \\ &= \sum_{k=1}^K -\frac{k}{2} \cdot \log(2\pi) - \frac{1}{2} \log(|\alpha^{-1} \cdot I|) - \frac{\alpha}{2} w_k^\top w_k = \\ &= -\frac{1}{2} \log(|\alpha^{-1} \cdot I|) + \sum_{k=1}^K -\frac{k}{2} \log(2\pi) - \frac{\alpha}{2} w_k^\top w_k \end{aligned}$$

, which is the log-prior.

The reason why we need the logarithm is because it performs a normalization to our values and scales them in a way that is more computer friendly. For example, very low probabilities are converted to manageable negative numbers. Also, the shift from the product " \prod " to the sum " \sum " can ~~scale~~ heavily potential small or big numbers to a manageable state.

(c) We want on to ~~compute~~ write a formula for the prior. From Bayes Rule we get:

$$P(w_1, \dots, w_k | \Phi, T, a) = \frac{P(T | \Phi, w_1, \dots, w_k) \cdot P(w_1, \dots, w_k | a)}{P(T | \Phi, a)} = \\ = \frac{\prod_{n=1}^N P(\vec{t}_n | \vec{\Phi}_n, w_1, \dots, w_k) \cdot P(w_1, \dots, w_k | a)}{\int p(T | \Phi, w_1, \dots, w_k) \cdot p(w_1, \dots, w_k | a) dw_1 \dots w_k}$$

③ In the final step all we did was a marginalization for w_1, \dots, w_k .

(d) From theory we know that the maximum a posteriori estimate is calculated via the sum of the log-likelihood and the log-prior, which we then solve in order to maximize W . So we have:

$$\operatorname{argmax}_w (\log(p(w_1, \dots, w_k | \Phi, T, a))) = \operatorname{argmax}$$

$$= \operatorname{argmax}_w \left[\prod_{n=1}^N P(\vec{t}_n | \vec{\Phi}_n, w_1, \dots, w_k) \cdot P(w_1, \dots, w_k | a) \right]$$

$$= \operatorname{argmax}_w \left[-\frac{1}{2} \cdot \log(|a^{-1}I|) + \sum_{k=1}^K \left(-\frac{k}{2} \log(2\pi) - \frac{\alpha}{2} w_k^\top w_k \right) + \sum_{n=1}^N \sum_{k=1}^K t_{nk} \cdot \log(y_k(\vec{\Phi}_n)) \right]$$

$$\text{remove constants} \quad \operatorname{argmax}_w \left[\sum_{k=1}^K \left(-\frac{\alpha}{2} w_k^\top w_k \right) + \sum_{n=1}^N \sum_{k=1}^K (t_{nk} \cdot \log(y_k(\vec{\Phi}_n))) \right] =$$

$$w^\top w = \|w\|_2^2 \quad = -\frac{\alpha}{2} \sum_{k=1}^K (\|w_k\|_2^2) + \sum_{n=1}^N \sum_{k=1}^K t_{nk} \cdot \log(y_k(\vec{\Phi}_n)) \quad (1)$$

(1) is identical to the expression of regularized logistic regression for K classes given by the exercise, but with opposite signs. Since we are trying to maximize (1) for W and the aim for the other expression is to minimize for W , then we can say that they are essentially solving the same problem.

$$\left(\sum_{n=1}^N \sum_{k=1}^K \left(t_{nk} \cdot \log(y_k(\vec{\phi}_n)) \right) - \frac{\alpha}{2} \cdot \sum_{k=1}^K (\|w_k\|_2^2) \right)$$

(e) To begin with, if we are quite confident that the weights should lie around 0, it means that our prior is strong and we should increase its impact in some way. This is done by the parameter "α", which is the precision. Since the weights lie around the mean, it means that we have low variance and therefore high precision. Thus, in the expression from (d) "α" will be large and will make the rightmost term more significant than the other terms. The reason that this actually makes the weights go to 0, is that the term of the prior is essentially the L2 norm of the weights. The L2 norm has a purpose to regularize the model by reducing the weights, therefore it will add a penalty term that will increase as they grow larger. Hence, by making this term more impactful, the norm will make sure to keep the weights as low as possible and close to the mean which is 0.

(f) The MLE is:

- $\underset{w}{\operatorname{argmax}} \left(\log \left(\prod_{n=1}^N p(\vec{t}_n | \vec{\phi}_n, w_1, \dots, w_K) \right) \right)$

The MAP as we saw in d is:

- $\underset{w}{\operatorname{argmax}} \left(\log \left(\prod_{n=1}^N p(\vec{t}_n | \vec{\phi}_n, w_1, \dots, w_K) \right) + \log(p(w_1, \dots, w_K | \alpha)) \right)$

We can see that if the prior term is 0 then we would have MLE = MAP. The uncertainty about the value of the weights described in the exercise, will make "α", the precision, to have a very low value, which will make the prior term (see below) to go to 0. The more uncertain we are, the more "α" decreases and thus, the more MAP approaches MLE.

$$\text{Prior term: } \underset{w}{\operatorname{argmax}} \log(p(w_1, \dots, w_K | \alpha)) = -\frac{\alpha}{2} \cdot \sum_{k=1}^K \|w_k\|_2^2$$

(g) We will derive the log-likelihood from (a).

$$\frac{\partial}{\partial w_j} \left(\sum_{n=1}^N \sum_{k=1}^K t_{nk} \cdot \log(y_k(\vec{\varphi}_n)) \right) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \cdot \frac{\partial}{\partial w_j} (\log(y_k(\vec{\varphi}_n))) = \\ = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \cdot \frac{1}{y_k(\vec{\varphi}_n)} \cdot (y_k(\vec{\varphi}_n)) \quad \textcircled{1}$$

$$\frac{\partial y_k}{\partial a_j} = y_k(\mathbb{I}_{kj} - y_j) \quad \text{(1)} \quad \text{and we have this from the exercise}$$

$$\text{We want to compute } \frac{\partial y_k}{\partial w_j} \xrightarrow[\text{rule?}]{\text{chain rule}} \sum_{p=1}^P \left(\frac{\partial y_k}{\partial a_p} \cdot \frac{\partial a_p}{\partial w_j} \right) \quad \text{(2)}$$

$$\frac{\partial a_p}{\partial w_j} = \frac{\partial}{\partial w_j} (w_p^\top \Phi) = \cancel{\Phi^\top} \delta p_j \quad \text{(3)} \quad \text{compute the other part}$$

$$\text{(2)} \xrightarrow[\text{(3)}]{=} \frac{\partial y_k}{\partial w_j} = \sum_{p=1}^P \left(y_k(\mathbb{I}_{kj} - y_j) \cdot \Phi^\top \delta p_j \right) = \quad \leftarrow \text{replace the derivatives}$$

$$= \sum_p^K \left(y_k \cdot \mathbb{I}_{kj} \cdot \Phi^\top \delta p_j - y_k \cdot y_j \cdot \Phi^\top \delta p_j \right) =$$

$$= \sum_p^K \left(y_k \cdot \mathbb{I}_{kj} \Phi^\top \delta p_j \right) - \sum_p^K \left(y_k \cdot y_j \cdot \Phi^\top \delta p_j \right) =$$

$$= y_k \cdot \mathbb{I}_{kj} \cdot \Phi^\top - y_k \cdot y_j \cdot \Phi^\top = \boxed{y_k(\vec{\varphi}_n) \cdot \Phi^\top (\mathbb{I}_{kj} - y_j) \quad \textcircled{2}}$$

$$\textcircled{2} \xrightarrow{\text{④}} \sum_{n=1}^N \sum_{k=1}^K t_{nk} \cdot \frac{1}{y_k(\vec{\varphi}_n)} \cdot y_k(\vec{\varphi}_n) \cdot \Phi^\top (\mathbb{I}_{kj} - y_j) = \quad \leftarrow \text{put it back in the initial relation and solve}$$

$$= \sum_{n=1}^N \sum_{k=1}^K \left(t_{nk} \cdot \Phi^\top (\mathbb{I}_{kj} - y_j) \right) =$$

$$= \sum_{n=1}^N \sum_{k=1}^K \left(t_{nk} \cdot \Phi^\top \mathbb{I}_{kj} - t_{nk} \cdot \Phi^\top y_j \right) =$$

$$= \sum_{n=1}^N \sum_{k=1}^K \left(t_{nk} \cdot \Phi^\top \mathbb{I}_{kj} \right) - \sum_{n=1}^N \sum_{k=1}^K \left(t_{nk} \cdot \Phi^\top y_j \right) = \quad \sum_{k=1}^K t_{nk} = 1$$

$$= \sum_{n=1}^N t_{nk} \cdot \Phi^\top - \sum_{n=1}^N \Phi^\top y_j \quad \text{⑤}$$

$$= \sum_{n=1}^N \Phi^\top (t_{nk} - y_j)$$

Thus, we have computed the derivative of the log-likelihood

(ii) Compute the derivative of the \log prior:

$$\frac{\partial}{\partial w_j} \left(\log(p(w_1, \dots, w_k | a)) \right) =$$

$$= \frac{\partial}{\partial w_j} \left[-\frac{1}{2} \left(\log(\alpha^{-1} I) + \sum_{k=1}^K k \cdot \log(2\pi + \alpha w_k^T w_k) \right) \right] =$$

$$= \frac{\partial}{\partial w_j} \left(\sum_{k=1}^K -\frac{\alpha}{2} w_k^T w_k \right) = -\frac{\alpha}{2} \sum_{k=1}^K \frac{\partial}{\partial w_j} (w_k^T w_k) =$$

$$= -\frac{\alpha}{2} \sum_{k=1}^K 2 w_k \cdot \delta_{kj} = -\frac{\alpha}{2} \cdot 2 w_j = -\alpha \cdot w_j$$

Compute the derivative of the \log posterior:

$$\frac{\partial}{\partial w_j} \left(\log \left(\prod_{u=1}^n p(t_u^* | \phi_u, w_1, \dots, w_k) \right) + \log(p(w_1, \dots, w_k | a)) \right) =$$

$$= \frac{\partial}{\partial w_j} \log \left(\prod_{u=1}^n p(t_u^* | \phi_u, w_1, \dots, w_k) \right) + \frac{\partial}{\partial w_j} \log(p(w_1, \dots, w_k | a)) =$$

substitute gradients $\sum_{u=1}^n \phi^T (t_{uk} - y_j(\phi_u)) - \alpha w_j$

Hence, we have computed the gradient for the log-prior and the log-posterior.

② (a) Logistic regression with linear features tries to find a straight line decision boundary that best separates the classes. By observing the dataset we can say that it is not a linearly separable dataset, which means that a straight line as a decision boundary would be prone to misclassifications. Thus, logistic regression with linear features won't be able to correctly classify the dataset.

(ii) Logistic regression with non-linear basis functions extends the capabilities of the basic logistic regression we saw previously. If we assign non-linear basis functions, logistic regression will be able to classify the dataset correctly. What we essentially want is an ellipse-like shape to maintain inside it one class and everything outside of it will belong to the other class. I believe that if we add quadratic basis functions, logistic regression will be able to "draw" this ellipsoid.

(iii) An MLP with 1-hidden layer as we can see in Bishop 4.17. can only "draw" straight lines. However, if we use a non-linear activation function then an MLP can approximate any decision boundary, which is exactly what makes it so powerful. The MLP will be able to make this classification effectively and correctly. A potential choice for an activation function is the sigmoid function (which obviously is non-linear).

(b) We consider the subset $\{(-1, 1), (1, 1), (-1, -1), (1, -1)\}$ and we will prove that its not possible to classify this using Naive Bayes.

Let's take the point $(1, 1)$ and classify it as C_1 . This means that:

$$\cdot P(C_1 | X_1=1, X_2=1) > P(C_2 | X_1=1, X_2=1) \quad (1)$$

We repeat this for the other 3 points as well:

$$\cdot P(C_1 | X_1=-1, X_2=-1) > P(C_2 | X_1=-1, X_2=-1) \quad (2)$$

$$\cdot P(C_1 | X_1=-1, X_2=1) < P(C_2 | X_1=-1, X_2=1) \quad (3)$$

$$\cdot P(C_1 | X_1=1, X_2=-1) < P(C_2 | X_1=1, X_2=-1) \quad (4)$$

We classified each point according to the shape, separating the classes "in diagonals" (s.t. each class is assigned 2 points in a non-linearly separable manner)

Then, we apply the Bayes rule as well as the Naive Bayes assumption to each relation.

$$(1) \Rightarrow p(x_1=1|C_1) \cdot p(x_2=1|C_1) \cdot p(C_1) > p(x_1=1|C_2) \cdot p(x_2=1|C_2) \cdot p(C_2)$$

$$\Leftrightarrow \frac{p(x_1=1|C_1) \cdot p(x_2=1|C_1) \cdot p(C_1)}{p(x_1=1|C_2) \cdot p(x_2=1|C_2) \cdot p(C_2)} > 1 \quad (5)$$

$$(2) \Rightarrow p(x_1=-1|C_1) \cdot p(x_2=-1|C_1) \cdot p(C_1) > p(x_1=-1|C_2) \cdot p(x_2=-1|C_2) \cdot p(C_2)$$

$$\Leftrightarrow \frac{p(x_1=-1|C_1) \cdot p(x_2=-1|C_1) \cdot p(C_1)}{p(x_1=-1|C_2) \cdot p(x_2=-1|C_2) \cdot p(C_2)} > 1 \quad (6)$$

$$(3) \Rightarrow p(x_1=-1|C_1) \cdot p(x_2=1|C_1) < p(x_1=-1|C_2) \cdot p(x_2=1|C_2) \cdot p(C_2)$$

$$\Leftrightarrow \frac{p(x_1=-1|C_2) \cdot p(x_2=1|C_2) \cdot p(C_2)}{p(x_1=-1|C_1) \cdot p(x_2=1|C_1) \cdot p(C_1)} > 1 \quad (7)$$

$$(4) \Rightarrow p(x_1=1|C_1) \cdot p(x_2=-1|C_1) < p(x_1=1|C_2) \cdot p(x_2=-1|C_2) \cdot p(C_2)$$

$$\Leftrightarrow \frac{p(x_1=1|C_2) \cdot p(x_2=-1|C_2) \cdot p(C_2)}{p(x_1=1|C_1) \cdot p(x_2=-1|C_1) \cdot p(C_1)} > 1 \quad (8)$$

Subsequently, we multiply (5) with (7) and (6) with (8) :

$$(5) \times (7) \Rightarrow \frac{p(x_1=1|C_1) \cdot p(x_2=1|C_1) \cdot p(C_1) \cdot p(x_1=-1|C_2) \cdot p(x_2=1|C_2) \cdot p(C_2)}{p(x_1=1|C_2) \cdot p(x_2=1|C_2) \cdot p(C_2) \cdot p(x_1=-1|C_1) \cdot p(x_2=1|C_1) \cdot p(C_1)} > 1$$

$$\Leftrightarrow \frac{p(x_1=1|C_1) \cdot p(x_1=-1|C_2)}{p(x_1=1|C_2) \cdot p(x_1=-1|C_1)} > 1 \quad (9)$$

$$(6) \times (8) \Rightarrow \frac{p(x_1=-1|C_1) \cdot p(x_2=-1|C_1) \cdot p(C_1) \cdot p(x_1=1|C_2) \cdot p(x_2=-1|C_2) \cdot p(C_2)}{p(x_1=-1|C_2) \cdot p(x_2=-1|C_2) \cdot p(C_2) \cdot p(x_1=1|C_1) \cdot p(x_2=1|C_1) \cdot p(C_1)} > 1$$

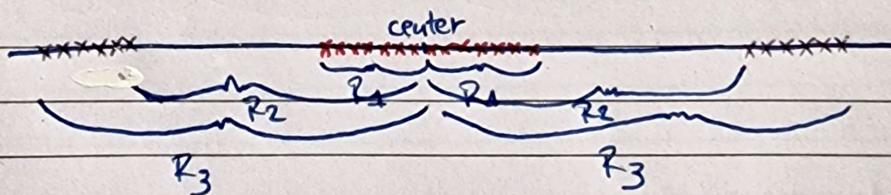
$$\Leftrightarrow \frac{p(x_1=-1|C_1) \cdot p(x_1=1|C_2)}{p(x_1=-1|C_2) \cdot p(x_1=1|C_1)} > 1 \quad (10)$$

We observe that (9) and (10) are reciprocals and it's impossible for them to be greater than 1. Therefore, we have arrived at a contradiction. The only assumption that we made throughout this whole proof was the Naive Bayes assumption. Thus, we have proved that Naive Bayes can not classify this dataset.

(c) The reasoning for this exercise will be different than (2b), because in contrast with that question, I believe that Naive Bayes can classify this dataset.

→ The dataset given in these ~~two~~ exercises causes 2 significant problems that prevent Naive Bayes from working. First, the classes are obviously non-linearly separable and second, there is intense correlation between x_1 and x_2 . I suggest that we switch the coordinates (x, y) to polar coordinates (r, θ) where it will be much simpler to separate linearly the classes. This transformation solves both issues, as in polar coordinates we will take advantage of the center of the data and the way they are circular around it (the orange dots are a uniform distribution, while the other class is called annulus in statistics), and we also solve the correlation issue, as the radius "r" tells us nothing about the angle " θ " and vice versa. Also, we choose to discard the feature θ altogether, as it won't contribute in any way to the classification and performing Naive Bayes in one dimension with only the radius will be much simpler. This scheme is depicted below:

Class 1: ●
Class 2: ●



With the above implementation, every point with radius up to R_2 will be classified as 'class 1', while all the rest can be classified as 'class 2'. In fact, this implementation can work for more than 2 classes by taking advantage of the other radii as well, as long as they maintain the same center and the same uniform formulation such as the classes that we have.

(d) The main difference between logistic regression with non-linear basis functions Φ , and multilayer perceptrons in terms of Φ , is that we manually choose static non-linear basis functions for logistic regression, while in the case of MLP's, the functions are ~~set through~~^{learned through} determined by the model weights that are updated in every epoch. At first the MLP chooses random weights that set random basis functions. In each epoch the weights are adjusted and then they adjust the basis functions. This is the main difference with logistic regression, as there we have to select the non-linear basis functions beforehand and they don't change during the process. In fact, this is what makes MLPs so attractive, as it is not always easy (rarely) to know in advance the proper basis functions.