

Mid-sem Exam (2017-2018, Sem-II)

Instructions:

- This question paper is worth a total of 100 marks.
- Total duration is 2.5 hours.
- Please write your name and roll number at the top of this sheet as well as on the provided answer copy.
- The **true-false** questions have to be answered on this question sheet itself, in the space provided. Please also submit the question sheet.
- To answer the other parts of the question paper, please use the provided answer copy.
- For rough work, please use pages at the back of answer copy. Extra sheets can be provided if needed.

Some distributions and their properties:

- For $x \in (0, 1)$, $\text{Beta}(x|a, b) = \frac{x^{a-1}(1-x)^{b-1}}{B(a, b)}$, where $B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$ and Γ denotes the gamma function s.t. $\Gamma(x) = (x-1)!$ for a positive integer x . Expectation of a Beta r.v.: $\mathbb{E}[x] = \frac{a}{a+b}$.
- For $x \in \{0, 1, 2, \dots\}$ (non-negative integers), $\text{Poisson}(x|\lambda) = \frac{\lambda^x \exp(-\lambda)}{x!}$ where λ is the rate parameter.
- For $x \in \mathbb{R}_+$, $\text{Gamma}(x|a, b) = \frac{b^a}{\Gamma(a)} x^{a-1} \exp(-bx)$ (shape and rate parameterization)
- For $x \in \mathbb{R}$, Univariate Gaussian: $\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\{-\frac{(x-\mu)^2}{2\sigma^2}\}$
- For $x \in \mathbb{R}^D$, D -dimensional Gaussian: $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\}$.
Trace-based representation: $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp\left\{-\frac{1}{2}\text{trace}[\boldsymbol{\Sigma}^{-1}\mathbf{S}]\right\}$, $\mathbf{S} = (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top$.
- For $\boldsymbol{\pi} = [\pi_1, \dots, \pi_K]$, s.t. $\sum_{k=1}^K \pi_k = 1$, $\text{Dirichlet}(\boldsymbol{\pi}|\alpha_1, \dots, \alpha_K) = \frac{1}{B(\alpha_1, \dots, \alpha_K)} \prod_{k=1}^K \pi_k^{\alpha_k-1}$ where $B(\alpha_1, \dots, \alpha_K) = \frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma(\sum_{k=1}^K \alpha_k)}$, and $\mathbb{E}[\pi_k] = \frac{\alpha_k}{\sum_{k=1}^K \alpha_k}$
- For $x_k \in \{0, N\}$ and $\sum_{k=1}^K x_k = N$, $\text{multinomial}(x_1, \dots, x_K|N, \boldsymbol{\pi}) = \frac{N!}{x_1! \dots x_K!} \pi_1^{x_1} \dots \pi_K^{x_K}$, where $\boldsymbol{\pi} = [\pi_1, \dots, \pi_K]$, s.t. $\sum_{k=1}^K \pi_k = 1$. The multinoulli is the same as multinomial with $N = 1$.

Some other useful results:

- If $\mathbf{x} = \mathbf{A}\mathbf{z} + \mathbf{b} + \epsilon$, $p(\mathbf{z}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$, $p(\epsilon) = \mathcal{N}(\epsilon|\mathbf{0}, \mathbf{L}^{-1})$ then $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{A}\mathbf{z} + \mathbf{b}, \mathbf{L}^{-1})$, $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^\top + \mathbf{L}^{-1})$, and $p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\Sigma}\{\mathbf{A}^\top \mathbf{L}(\mathbf{x} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma})$, where $\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^\top \mathbf{L} \mathbf{A})^{-1}$.
- Marginal and conditional distributions for Gaussians: $p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa})$, $p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b})$ where $\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Lambda}_{aa}^{-1} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba}$, $\boldsymbol{\mu}_{a|b} = \boldsymbol{\Sigma}_{a|b}\{\boldsymbol{\Lambda}_{aa}\boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b)\} = \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1}\boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) = \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b)$, where symbols have their usual meaning. :)
- $\frac{\partial}{\partial \boldsymbol{\mu}}[\boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu}] = [\mathbf{A} + \mathbf{A}^\top]\boldsymbol{\mu}$, $\frac{\partial}{\partial \mathbf{A}} \log |\mathbf{A}| = \mathbf{A}^{-\top}$, $\frac{\partial}{\partial \mathbf{A}} \text{trace}[\mathbf{A}\mathbf{B}] = \mathbf{B}^\top$
- For a random variable vector \mathbf{x} , $\mathbb{E}[\mathbf{x}\mathbf{x}^\top] = \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{x}]^\top + \text{cov}[\mathbf{x}]$

Questions begin from the next page

1 True-False Problems (10 questions, 10 marks)

1. [T] A likelihood distribution $p(x|\theta)$ does not necessarily sum/integrate to 1 over all possible values of the parameter θ .
2. [F] Assuming continuous-valued model parameters, the posterior predictive distribution can be seen as a uniformly weighted average of infinite many point predictive distributions.
3. [T] The log-likelihood $\log p(\mathbf{X}|\Theta)$ can never change in the E step of the EM algorithm.
4. [T] Marginal posterior and conditional posterior are defined in a lower dimensional parameter space as compared to a joint posterior.
5. [F] The computational cost of fully Bayesian inference is independent of the amount of data and only depends on the number of parameters and latent variables.
6. [F] When using the Bayes rule to find the expression for the posterior, the denominator in the Bayes rule expression (i.e., the marginal likelihood) can be ignored.
7. [T] In GP regression/classification, the GP prior does not depend on the data being modeled.
8. [F] Generative classification in general requires fewer parameters to estimate as compared to discriminative classification.
9. [F] Choosing the best model m based on the comparing the posterior probabilities of different models is equivalent to choosing the model with the largest marginal likelihood $p(\mathbf{X}|m)$.
10. [T] For probabilistic PCA with $p(\mathbf{x}_n|\mathbf{z}_n) = \mathcal{N}(\mathbf{W}\mathbf{z}_n, \sigma^2\mathbf{I}_D)$ with $p(\mathbf{z}_n) = \mathcal{N}(0, \mathbf{I}_K)$, the “spread” of latent variable’s posterior $p(\mathbf{z}_n|\mathbf{x}_n)$ does not depend on the input \mathbf{x}_n .

2 Short Answer Problems (10 questions, 30 marks)

1. Rank the three parameter estimation approaches, namely maximum-likelihood estimation (MLE), maximum-a-posteriori (MAP) estimation, and fully Bayesian inference, based on their robustness against overfitting, and give a brief justification for your ranking.
Answer: The order will be (maximum prone to least prone)- MLE > MAP > Fully Bayesian inference. MLE: point estimate, no prior/regularization. MAP: Point estimate but incorporates prior/regularizer. Fully Bayesian inference: Full posterior distribution and posterior averaging at prediction time.
2. Suppose a probability distribution $p(\mathbf{x}|\theta)$ has a form proportional to $\exp(-\mathbf{x}^\top \mathbf{A} \mathbf{x} - \mathbf{b}^\top \mathbf{x})$, where $\mathbf{x} \in \mathbb{R}^D$. Identify this distribution (i.e., its name) and write down the expressions for its parameters θ .
Answer: One quadratic and one linear term in x , so it must be a Gaussian $\mathcal{N}(\mu, \Sigma)$. Comparing with forms, we get $\Sigma = (2\mathbf{A})^{-1}$ and $\mu = -(2\mathbf{A})^{-1}\mathbf{b}$
3. Intuitively, which of the two distributions - plug-in predictive distribution and posterior predictive distribution - is expected to have a larger variance at a test point? Briefly justify your answer.
Answer: Posterior predictive is expected to have a larger variance since it uses the entire posterior of the parameters to make predictive (thus taking into account their uncertainty while making prediction), rather than relying on a single value of the parameter to make prediction (as done by plug-in predictive)
4. Consider data with N training examples $\{\mathbf{x}_n, y_n\}_{n=1}^N$, with $\mathbf{x}_n \in \mathbb{R}^D$, and $D \gg N$. Suppose we have learned two models - GP regression and Bayesian linear regression. For calculating the mean of the posterior predictive, which of the two models will be faster, in terms of big-O notation, and why?
Answer: GP takes $O(N)$ time to compute the mean of the posterior predictive (predictive mean of the form $\sum_{n=1}^N \alpha_n k(\mathbf{x}_n, \mathbf{x}_*)$), whereas Bayesian linear regression takes $O(D)$ time (predictive mean of the form $\mu_N^\top \mathbf{x}_*$ where $\mu_N \in \mathbb{R}^D$ is the mean of the posterior). If $D \gg N$, GP prediction will be faster.

5. Give two examples where Laplace approximation may not give a good approximation for a probability distribution. You may draw figures to illustrate.

Answer: Multimodal, asymmetric shaped distributions, distributions with only positive support, etc.

6. Consider two Gaussians $\mathcal{N}(x|\mu_1, \Sigma_1)$ and $\mathcal{N}(x|\mu_2, \Sigma_2)$. Suppose we use the following rule to assign a binary label y to a point x : if $\mathcal{N}(x|\mu_1, \Sigma_1) \geq \mathcal{N}(x|\mu_2, \Sigma_2)$ then $y = +1$, otherwise $y = -1$.

Under what condition we can express the above rule as $y = \text{sign}[f(x)]$ where $f(x)$ is a *linear* function of x (i.e., of the form $w^\top x$), and what will be the expression for $f(x)$ in such a case?

Answer: Given the condition, we must have $\log \mathcal{N}(x|\mu_1, \Sigma_1) \geq \log \mathcal{N}(x|\mu_2, \Sigma_2)$. Expanding, we get $-0.5 \log |\Sigma_1^{-1}| - (x - \mu_1)^\top \Sigma_1^{-1} (x - \mu_1) > -0.5 \log |\Sigma_2^{-1}| - (x - \mu_2)^\top \Sigma_2^{-1} (x - \mu_2)$. The quadratic terms in x will vanish if $\Sigma_1 = \Sigma_2 = \Sigma$ (say), and we will be left with $w^\top x + b \geq 0$, where $w = \Sigma^{-1}(\mu_1 - \mu_2)$ and $b = 0.5(\mu_2^\top \Sigma^{-1} \mu_2 - \mu_1^\top \Sigma^{-1} \mu_1)$. Thus $f(x) = w^\top x + b$.

7. Briefly explain why marginal likelihood can be seen as a special case of a posterior predictive.

Answer: Posterior predictive integrates out the likelihood over the posterior ($p(y_*|\mathbf{y}) = \int p(y_*|\theta)p(\theta|\mathbf{y})d\theta$) and marginal likelihood does it over the prior ($p(y_*) = \int p(y_*|\theta)p(\theta)d\theta$). Therefore, marginal likelihood can be thought of as the posterior predictive for y_* with $N = 0$, i.e., no “training” data \mathbf{y}

8. When using EM for probabilistic PCA, can the dimensionality of the latent space be more than the data dimensionality? If yes, why? If no, why not?

Answer: Yes. There is nothing in the PPCA model definition or its EM algorithm, that requires $K \leq D$ (unlike standard PCA which uses eigendecomposition and can therefore can't have $K > D$).

9. Briefly explain why, for EM, the M step of finding the optimal Θ is a convex problem, if the joint distribution of data and latent variables given parameters, i.e., $p(\mathbf{x}_n, \mathbf{z}_n|\Theta)$ is in exponential family.

Answer: Note that the joint distribution will be of the form $p(\mathbf{x}_n, \mathbf{z}_n|\Theta) = h(\mathbf{x}_n, \mathbf{z}_n) \exp(\Theta^\top \phi(\mathbf{x}_n, \mathbf{z}_n)) - A(\Theta)$. The M step maximizes the expectation of the log of this quantity (summed over all observations). Thus (ignoring the expectation) the objective will be of the form $\Theta^\top \sum_{n=1}^N \phi(\mathbf{x}_n, \mathbf{z}_n) - NA(\Theta)$. Since $A(\Theta)$ is convex (and its negative concave) and the first term is linear, the overall objective is a concave function of Θ (and if we *minimize* it then the objective will be its negative, which is convex).

10. Consider a latent variable model where the prior over the latent variables is of the form $p(\mathbf{z}_n|\mathbf{z}_{n-1})$, i.e., the distribution of \mathbf{z}_n depends on the *value* of the latent variable of the previous observation. Assuming \mathbf{z}_n to be discrete with K possibilities, which distribution would you use to model $p(\mathbf{z}_n|\mathbf{z}_{n-1})$ and how many parameters would be needed to model this $p(\mathbf{z}_n|\mathbf{z}_{n-1})$?

Answer: Note that, since \mathbf{z}_n is discrete with K possibilities, for each value (say ℓ) of \mathbf{z}_{n-1} we can model $p(\mathbf{z}_n|\mathbf{z}_{n-1} = \ell)$ using a multinoulli with K dimensional probability vector $\boldsymbol{\pi}_\ell$. We will need K such probability vectors (one for each value of \mathbf{z}_{n-1}) and therefore will need total K^2 parameters.

3 Medium-Length Answer Problems (5 questions, 40 marks)

1. The KL divergence between two discrete distributions $p(x)$ and $q(x)$ with their support being the set of the non-negative integers is defined as $KL(p||q) = \sum_{x=0}^{\infty} p(x) \log \frac{p(x)}{q(x)}$. If $q(x)$ is the Poisson distribution, i.e., $q(x) = \frac{\lambda^x \exp(-\lambda)}{x!}$ then what value of λ will minimize the KL divergence?

Answer: Minimizing the KL w.r.t. $q(x)$, which in turn depends on λ will give

$$\hat{\lambda} = \arg \min_{\lambda} - \sum_{x=0}^{\infty} p(x) \log q(x) = \arg \min_{\lambda} - \sum_{x=0}^{\infty} p(x) [x \log \lambda - \lambda]$$

Differentiating w.r.t. λ and setting it to zero gives $\hat{\lambda} = \sum_{x=0}^{\infty} xp(x) = \mathbb{E}[x]$. (note: you can check that the second derivative will be positive, so it must be the minimizer)

2. Express the Bernoulli $p(x|\theta) = \theta^x(1-\theta)^{1-x}$ in the exponential family form. Write down the expressions for its natural parameters, the sufficient statistics, and the log-partition function. Also show that the first derivative of its log-partition function is equal to the expected value of its sufficient statistics.

Answer: Rewriting $p(x|\theta)$ as $\exp(\log p(x|\theta))$ and comparing with the form of an exponential family distribution $p(x|\eta) = h(x) \exp(\eta^\top \phi(x) - A(\eta))$, we get $\eta = \log \frac{\theta}{1-\theta}$ as the natural parameters, $\phi(x) = x$ as the sufficient statistics, and $A(\eta) = -\log(1-\theta) = \log(1 + \exp(\eta))$ as the log-partition function. The first derivative of $A(\eta) = \frac{\exp(\eta)}{1+\exp(\eta)} = \theta = \mathbb{E}[x]$

3. Consider a Gaussian Mixture Model (GMM) with K components where the posterior probabilities of cluster memberships for a data point \mathbf{x}_n are defined as

$$p(z_n = k|\mathbf{x}_n) = \gamma_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}|\mu_k, \sigma^2 \mathbf{I})}{\sum_{\ell=1}^K \pi_\ell \mathcal{N}(\mathbf{x}|\mu_\ell, \sigma^2 \mathbf{I})}$$

where π_k is the mixing proportion for k^{th} Gaussian, μ_k is its mean, and $\sigma^2 \mathbf{I}$ is its covariance matrix (assumed the same for all the Gaussians). Show that as σ^2 tends to zero, the EM algorithm for this GMM reduces to a K -means style algorithm. Also briefly summarize the steps of this resulting EM algorithm, e.g., what would each step compute, writing down the necessary expressions.

Answer: The probability of \mathbf{x}_n belonging to cluster k computed in E step is

$$\gamma_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}|\mu_k, \sigma^2 \mathbf{I})}{\sum_{\ell=1}^K \pi_\ell \mathcal{N}(\mathbf{x}|\mu_\ell, \sigma^2 \mathbf{I})} = \frac{\pi_k \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{\|\mathbf{x}_n - \mu_k\|^2}{2\sigma^2})}{\sum_{\ell=1}^K \pi_\ell \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{\|\mathbf{x}_n - \mu_\ell\|^2}{2\sigma^2})}$$

Now, as σ^2 tends to zero, out the K Gaussians in the denominator, the sum will be dominated by the Gaussian for which $\|\mathbf{x}_n - \mu_\ell\|^2$ is the smallest. In that case, the cluster membership probability γ_{nk} will be close to 1 for that Gaussian and close to 0 for all other Gaussians. Thus the vector $\gamma = [\gamma_{n1}, \dots, \gamma_{nK}]$ will be like a one-hot vector and GMM will reduce to a hard clustering algorithm. The EM algorithm would alternate between the following steps:

- E step: For each $n = 1, \dots, N$, suppose $k_* = \arg \min_\ell \|\mathbf{x}_n - \mu_\ell\|^2$. Set $z_n = k_*$
- M step: For each $k = 1, \dots, K$, set μ_k to be the mean of all points assigned to cluster k

Although you can also estimate π_k as $\pi_k = \frac{\sum_{n: z_n=k} \mathbf{x}_n}{N}$, the EM algorithm doesn't require π_k .

Note: This technique where we take a probabilistic model and let the variance in the likelihood term go to zero is known as “small-variance asymptotics” (SVA), and is a popular way to convert a probabilistic model into a similar non-probabilistic model that has a faster inference. It is not limited to Gaussian likelihoods where variance has a natural interpretation but works for any exponential family distribution. It is also used as a fast point estimation algo. for nonparametric Bayesian models.

4. Consider a mixture model $p(\mathbf{x}) = \sum_{k=1}^K \pi_k p(\mathbf{x}|k)$ where $p(\mathbf{x}|k)$ denotes the k^{th} component distribution of this mixture model and π_k denotes its mixing proportion. Assume each observation \mathbf{x} to consist of two parts, i.e., $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2]$. Show that the conditional distribution $p(\mathbf{x}_1|\mathbf{x}_2)$ can also be written in form of a mixture distribution and derive the expressions for the mixing proportions and the component distributions for this mixture distribution. (Hint: Note that you can write $p(\mathbf{x}_1|\mathbf{x}_2)$ as the ratio of a joint distribution and a marginal distribution)

Answer: Note that $p(\mathbf{x}) = p(\mathbf{x}_1, \mathbf{x}_2) = \sum_{k=1}^K \pi_k p(\mathbf{x}_1, \mathbf{x}_2|k)$. Therefore

$$p(\mathbf{x}_1|\mathbf{x}_2) = \frac{p(\mathbf{x}_1, \mathbf{x}_2)}{p(\mathbf{x}_2)} = \frac{\sum_{k=1}^K \pi_k p(\mathbf{x}_1, \mathbf{x}_2|k)}{\int \sum_{\ell=1}^K \pi_\ell p(\mathbf{x}_1, \mathbf{x}_2|\ell) d\mathbf{x}_1} = \frac{\sum_{k=1}^K \pi_k p(\mathbf{x}_1, \mathbf{x}_2|k)}{\sum_{\ell=1}^K \pi_\ell p(\mathbf{x}_2|\ell)}$$

Finally, we can always write $p(\mathbf{x}_1, \mathbf{x}_2|k)$ as $p(\mathbf{x}_1, \mathbf{x}_2|k) = p(\mathbf{x}_1|\mathbf{x}_2, k)p(\mathbf{x}_2|k)$, we have

$$p(\mathbf{x}_1|\mathbf{x}_2) = \frac{\sum_{k=1}^K \pi_k p(\mathbf{x}_1|\mathbf{x}_2, k)p(\mathbf{x}_2|k)}{\sum_{\ell=1}^K \pi_\ell p(\mathbf{x}_2|\ell)} = \sum_{k=1}^K \eta_k(\mathbf{x}_2) p(\mathbf{x}_1|\mathbf{x}_2, k)$$

where $\eta_k(\mathbf{x}_2) = \frac{\pi_k p(\mathbf{x}_2|k)}{\sum_{\ell=1}^K \pi_\ell p(\mathbf{x}_2|\ell)}$ denotes the k^{th} mixing proportion and $p(\mathbf{x}_1|\mathbf{x}_2, k)$ the k^{th} component distribution. Note that the mixing proportion as well as the component distribution depend on \mathbf{x}_2 .

5. Consider a Gaussian Process (GP) prior $p(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})$ where $\mathbf{f} = [f_1, \dots, f_N]$ is the vector of scores of a function f on each of the N inputs $\mathbf{x}_1, \dots, \mathbf{x}_N$ and \mathbf{K} is the $N \times N$ kernel/covariance matrix with $K_{nm} = \kappa(\mathbf{x}_n, \mathbf{x}_m)$ being the pairwise similarity between inputs \mathbf{x}_n and \mathbf{x}_m .

Assume the following likelihood model for each response $p(y_n|f_n) = \mathcal{N}(y_n|f_n, \beta^{-1})$. Denote $\mathbf{y} = [y_1, \dots, y_N]$ as the vector of all the responses and derive the expression for the posterior distribution $p(\mathbf{f}|\mathbf{y})$. Feel free to use properties of Gaussian distributions (don't derive from scratch).

Now consider another case when each response $y_n \in \{0, 1\}$. What would be the an appropriate likelihood model in this case? Suggest a method for derive the posterior $p(\mathbf{f}|\mathbf{y})$ in this case. You don't need to give the entire derivation – a brief outline of how you would do it in this case is fine.

Answer: The prior $p(\mathbf{f}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})$ and the likelihood $p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I})$. Since both are Gaussian, the posterior $p(\mathbf{f}|\mathbf{y})$ must also be Gaussian and is given by (using reverse conditioning)

$$p(\mathbf{f}|\mathbf{y}) \propto \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}) \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I}) = \mathcal{N}(\mathbf{f}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

where $\boldsymbol{\Sigma} = (\mathbf{K}^{-1} + \frac{1}{\sigma^2} \mathbf{I})^{-1}$ and $\boldsymbol{\mu} = \boldsymbol{\Sigma}[\frac{1}{\sigma^2} \mathbf{y}] = (\sigma^2 \mathbf{K}^{-1} + \mathbf{I}_N)^{-1} \mathbf{y}$.

For the case when y_n is binary, closed form posterior is not available. Since the problem is similar to logistic regression (though with GP prior), one way to approximate the posterior would be Laplace approximation. To do so, we can first find the MAP estimate of \mathbf{f} by maximizing the log-posterior $p(\mathbf{f}|\mathbf{y}) \propto \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}) \prod_{n=1}^N \text{Bernoulli}(y_n|\sigma(f_n))$. We can also take compute the Hessian of the log-posterior. The MAP and Hessian can be used as the mean and covariance matrix of the Gaussian used in the Laplace approximation of the posterior $p(\mathbf{f}|\mathbf{y})$.

4 Long Answer Problem (1 question, 20 marks)

Consider a regression problem setting where we don't just want to learn a single regression model on a single dataset but instead wish to learn T regression models on T datasets $\{\mathcal{D}_t\}_{t=1}^T$. Let's call each regression problem a "task", so there are a total of T tasks to be learned. Assume that each dataset $\mathcal{D}_t = (\mathbf{X}^{(t)}, \mathbf{y}^{(t)})$ consists of N_t training examples $(\mathbf{x}_n^{(t)}, y_n^{(t)})_{n=1}^{N_t}$, where $\mathbf{x}_n^{(t)} \in \mathbb{R}^D$ and $y_n^{(t)} \in \mathbb{R}$.

Assume each dataset/task \mathcal{D}_t to be modeled by a linear regression model with weight vector $\mathbf{w}_t \in \mathbb{R}^D$. For simplicity, assume the Gaussian noise precision to be known and the same across all tasks. Our goal will be to learn the T regression weight vectors $\mathbf{w}_1, \dots, \mathbf{w}_T$. However, instead of learning each weight vector \mathbf{w}_t *independently* only using \mathcal{D}_t , we wish to learn all the weight vectors in a "collaborative" fashion, so that these tasks can "share information" with each other (which may be especially helpful if the amount of training data for each task is very small).

To do this, let us assume that each \mathbf{w}_t can be written as a linear combination of K vectors $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_K]$, i.e., $\mathbf{w}_t = \mathbf{B} \mathbf{z}_t = \sum_{k=1}^K z_{tk} \mathbf{b}_k$ where each $\mathbf{b}_k \in \mathbb{R}^D$, $k = 1, \dots, K$, and $\mathbf{z}_t \in \mathbb{R}^K$ are the weights of this linear combination for task t . Note that the $D \times K$ matrix $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_K]$ is shared by all the tasks (which helps sharing information across all the tasks) and $\mathbf{z}_t \in \mathbb{R}^K$ is specific to task t . As you can see, since $\mathbf{w}_t = \mathbf{B} \mathbf{z}_t$, learning $\mathbf{w}_1, \dots, \mathbf{w}_T$ now reduces to learning \mathbf{B} and $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_T]$.

For simplicity, assume zero mean, identity covariance Gaussian priors on each \mathbf{b}_k , $k = 1, \dots, K$ and each \mathbf{z}_t , $t = 1, \dots, T$, so $p(\mathbf{b}_k) = \mathcal{N}(\mathbf{b}_k|0, \mathbf{I}_D)$ and $p(\mathbf{z}_t) = \mathcal{N}(\mathbf{z}_t|0, \mathbf{I}_K)$.

You can think of this model as a multiparameter model with unknowns \mathbf{B} and \mathbf{Z} (treat all other hyperparameters as known), and can break the inference problem into inferring the conditional posteriors of each of the unknowns, i.e., $p(\mathbf{b}_k|\mathbf{B}_{-k}, \mathbf{Z}, \{\mathcal{D}_t\}_{t=1}^T)$, $k = 1, \dots, K$, and $p(\mathbf{z}_t|\mathbf{B}, \mathbf{Z}_{-t}, \{\mathcal{D}_t\}_{t=1}^T)$, $t = 1, \dots, T$. Here \mathbf{B}_{-k} denotes \mathbf{B} with \mathbf{b}_k excluded. Likewise, \mathbf{Z}_{-t} denotes \mathbf{Z} with \mathbf{z}_t excluded.

- Give the expressions for the conditional posteriors $p(\mathbf{b}_k|\mathbf{B}_{-k}, \mathbf{Z}, \{\mathcal{D}_t\}_{t=1}^T)$ and $p(\mathbf{z}_t|\mathbf{B}, \mathbf{Z}_{-t}, \{\mathcal{D}_t\}_{t=1}^T)$.
- Also, using the obtained expression/form of $p(\mathbf{z}_t|\mathbf{B}, \mathbf{Z}_{-t}, \{\mathcal{D}_t\}_{t=1}^T)$, explain why inferring the conditional posterior for each \mathbf{z}_t is equivalent to inferring the posterior of a Bayesian linear regression model using only \mathcal{D}_t , but with the inputs $\mathbf{x}_n^{(t)}$ projected to a K dimensional space via a projection matrix.

Hint: With the parametrization $\mathbf{w}_t = \mathbf{B}\mathbf{z}_t$ assumed in this problem, given \mathbf{B}_{-k} and \mathbf{Z} , one can think of \mathbf{b}_k as the weight vector of an appropriately defined Bayesian linear regression model. Likewise, given \mathbf{B} and \mathbf{Z}_{-t} , one can think of \mathbf{z}_t as the weight vector of an appropriately defined Bayesian linear regression problem.

Answer: We will be using an alternating inference scheme in which we will infer one unknown at a time, keeping all other fixed, and cycling through until we converge.

- Let's first consider inferring each \mathbf{b}_k . The likelihood model for task t can be written as $y_n^{(t)} = \mathbf{w}_t^\top \mathbf{x}_n^{(t)} + \epsilon_n^{(t)}$. Since $\mathbf{w}_t = \sum_{k=1}^K z_{tk} \mathbf{b}_k$, the likelihood can be written as

$$y_n^{(t)} = \mathbf{b}_k^\top (z_{tk} \mathbf{x}_n^{(t)}) + \sum_{\ell \neq k} \mathbf{b}_\ell^\top (z_{t\ell} \mathbf{x}_n) + \epsilon_n^{(t)}$$

Given \mathbf{Z} and \mathbf{B}_{-k} the summation above is a constant we can rewrite the above as

$$y_n^{(t)} - \sum_{\ell \neq k} \mathbf{b}_\ell^\top (z_{t\ell} \mathbf{x}_n) = \mathbf{b}_k^\top (z_{tk} \mathbf{x}_n^{(t)}) + \epsilon_n^{(t)}$$

Note that the above is a standard probabilistic regression model with a “modified” input $\hat{\mathbf{x}}_n^{(t)} = z_{tk} \mathbf{x}_n^{(t)}$ and a “modified” response $\hat{y}_n^{(t)} = y_n^{(t)} - \sum_{\ell \neq k} \mathbf{b}_\ell^\top (z_{t\ell} \mathbf{x}_n)$. We can therefore take the “modified” data $(\hat{\mathbf{x}}_n^{(t)}, \hat{y}_n^{(t)})$ from all the tasks $t = 1, \dots, T$, and learn a Bayesian linear regression model. The conditional posterior $p(\mathbf{b}_k | \mathbf{B}_{-k}, \mathbf{Z}, \{\mathcal{D}_t\}_{t=1}^T)$ will have a closed form since the prior on \mathbf{b}_k is Gaussian and the conditional likelihood for the modified version of the data is still Gaussian (note the equation above).

Also note that \mathbf{b}_k is shared by all the tasks, which is why the conditional posterior for \mathbf{b}_k will indeed depend on all the data from all the tasks (as can be seen from the likelihood model above).

- Now, let's consider inferring \mathbf{z}_t . We can again rewrite the likelihood model as

$$y_n^{(t)} = (\mathbf{B}\mathbf{z}_t)^\top \mathbf{x}_n^{(t)} + \epsilon_n^{(t)} = \mathbf{z}_t^\top (\mathbf{B}^\top \mathbf{x}_n^{(t)}) + \epsilon_n^{(t)}$$

Given \mathbf{Z}_{-t} and \mathbf{B} , the above is again a probabilistic linear regression model with weight vector \mathbf{z}_t and inputs $\mathbf{B}^\top \mathbf{x}_n^{(t)}$. Note that the above inputs are a K dimensional projection of the original inputs $\mathbf{x}_n^{(t)}$ by a projection matrix defined by the $D \times K$ matrix \mathbf{B} . The conditional posterior $p(\mathbf{z}_t | \mathbf{B}, \mathbf{Z}_{-t}, \{\mathcal{D}_t\}_{t=1}^T)$ will have a closed form since the prior on \mathbf{z}_t is Gaussian and the conditional likelihood for the modified version of the data is still Gaussian (note the equation above).

Note that since \mathbf{z}_t is specific to task t only, the conditional posterior for \mathbf{z}_t actually only depends on the data from task t , so the correct notation will be $p(\mathbf{z}_t | \mathbf{B}, \mathcal{D}_t)$; in fact, it also doesn't depend on \mathbf{Z}_{-t} (as can be seen from the likelihood model above).