

Paper 3F8: Inference

Solutions to Example Sheet 2: Bayesian Linear Regression,
Classification, Dimensionality Reduction and Clustering

Bayesian linear regression

1. Bayesian Linear Regression

A single data point $\{x, y\}$ is fit using Bayesian linear regression. The output y is assumed to be generated from the input x according to a linear relationship that is corrupted by Gaussian noise $y = mx + c + \epsilon$. The noise ϵ is mean 0 and variance 1 so $p(y|m, c, x) = \mathcal{N}(y; mx + c, 1)$. Gaussian priors are placed on the slope m and intercept c with zero mean and unit variance, that is $p(m) = \mathcal{N}(m; 0, 1)$ and $p(c) = \mathcal{N}(c; 0, 1)$.

- (a) Compute the posterior probability of the slope and intercept given the data point, that is $p(m, c|x, y)$.

$$a) p(m, c | y, x) = \frac{p(y|m, c, x) p(m)p(c)}{p(y|x)}$$

$$\propto p(y|m, c, x) p(m)p(c)$$

$= N\left(\begin{bmatrix} m \\ c \end{bmatrix} ; \underline{\mu}_{\text{post}}, \underline{\Sigma}_{\text{post}}\right)$

tactic ← product of Gaussian density in M&C

① Substitute in for these densities and ② identify mean & covariance by comparing coefficients

② ← posterior must be Gaussian: just need to compute mean & covariance

$$\begin{aligned} ① p(m, c | y, x, \sigma_y^2) &\propto \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \frac{(y - mx - c)^2}{\sigma_y^2}} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \frac{m^2}{\sigma_y^2}} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \frac{c^2}{\sigma_y^2}} \\ &\propto e^{-\frac{1}{2} \frac{y^2}{\sigma_y^2} - \frac{1}{2} \frac{m^2}{\sigma_y^2} - \frac{1}{2} \frac{c^2}{\sigma_y^2} - \frac{1}{2} \frac{2mx + 2c}{\sigma_y^2}} \\ &\propto e^{-\frac{1}{2} \left(\begin{bmatrix} m \\ c \end{bmatrix} - \underline{\mu}_{\text{post}} \right)^T \left(\underline{\Sigma}_{\text{post}}^{-1} \left(\begin{bmatrix} m \\ c \end{bmatrix} - \underline{\mu}_{\text{post}} \right) \right)} \\ &\propto e^{-\frac{1}{2} \left[\begin{bmatrix} m \\ c \end{bmatrix}^T \underline{\Sigma}_{\text{post}}^{-1} \begin{bmatrix} m \\ c \end{bmatrix} + \begin{bmatrix} m \\ c \end{bmatrix}^T \underline{\Sigma}_{\text{post}}^{-1} \underline{\mu}_{\text{post}} \right]} \\ &\Rightarrow \underline{\Sigma}_{\text{post}} = \begin{bmatrix} x^2 + 1 & x \\ x & 2 \end{bmatrix}^{-1} = \frac{1}{(x^2 + 1)(2) - x^2} \begin{pmatrix} 2 & -x \\ -x & x^2 + 1 \end{pmatrix} \\ &= \frac{1}{x^2 + 2} \begin{pmatrix} 2 & -x \\ -x & x^2 + 1 \end{pmatrix} \end{aligned}$$

inverse covariance times mean

inverse covariance

$$\& \text{since } \left(\underline{\Sigma}_{\text{post}} \right)^{-1} \underline{\mu}_{\text{post}} = \begin{bmatrix} yx \\ y \end{bmatrix}$$

$$\underline{\mu}_{\text{post}} = \frac{1}{x^2 + 2} \begin{pmatrix} 2 & -x \\ -x & x^2 + 1 \end{pmatrix} \begin{bmatrix} yx \\ y \end{bmatrix}$$

$$= \frac{1}{x^2 + 2} \begin{pmatrix} 2yx - yx \\ -yx^2 + (x^2 + 1)y \end{pmatrix} = \frac{1}{x^2 + 2} \begin{pmatrix} yx \\ y \end{pmatrix}$$

- (b) Show that the posterior derived in part a is consistent with the expressions for Bayesian linear regression given in lectures.

b) Check these expressions match with those in lectures : $\underline{y} = \underline{\tilde{x}} \underline{w} + \underline{\varepsilon}$ $\underline{\varepsilon} \sim N(0, \sigma^2 \underline{\underline{I}})$

$$p(\underline{w} | \text{Data}, \sigma^2, \lambda) = N(\underline{w}; \underline{\mu}_{\text{WID}}, \underline{\Sigma}_{\text{WID}}) \quad \underline{w} \sim N(0, \underline{\underline{x}}^\top \underline{\underline{x}})$$

$$\text{where } \underline{\Sigma}_{\text{WID}} = \left(\underline{\underline{\lambda}} + \frac{1}{\sigma^2} \underline{\tilde{x}}^\top \underline{\tilde{x}} \right)^{-1} \quad \underline{\mu}_{\text{WID}} = \underline{\Sigma}_{\text{WID}} \frac{1}{\sigma^2} \underline{\tilde{x}}^\top \underline{y}$$

in our case $\sigma^2 = \lambda = 1$ (prior & observation noise are unit variance)

$$\underline{\tilde{x}} = [x, 1]$$

$$\underline{w} = \begin{bmatrix} m \\ c \end{bmatrix}$$

$$\Rightarrow \underline{\Sigma}_{\text{WID}} = \left(\underline{\underline{\lambda}} + \begin{bmatrix} x \\ 1 \end{bmatrix} \begin{bmatrix} x & 1 \end{bmatrix} \right)^{-1} = \begin{pmatrix} x^2 + 1 & x \\ x & 2 \end{pmatrix}^{-1} \quad \checkmark$$

$$\underline{\mu}_{\text{WID}} = \begin{pmatrix} x^2 + 1 & x \\ x & 2 \end{pmatrix}^{-1} \begin{bmatrix} x \\ 1 \end{bmatrix} \underline{y} \quad \checkmark$$

(c) Compute the posterior in the following three cases and provide explanations for why the posteriors take the form that they do.

- i. $x = 0$ and $y = 0$
- ii. $x = 1$ and $y = 0$
- iii. $x = 100$ and $y = 100$

c) plug the following settings $\begin{matrix} x=0 \\ y=0 \end{matrix}$ $\begin{matrix} x=1 \\ y=1 \end{matrix}$ $\begin{matrix} x=100 \\ y=100 \end{matrix}$

into :

$$\underline{\Sigma}^{\text{post}} = \frac{1}{x^2+2} \begin{bmatrix} 2 & -x \\ -x & x^2+1 \end{bmatrix}$$

$$\underline{\mu}^{\text{post}} = \frac{1}{x^2+2} \begin{bmatrix} y \\ x \end{bmatrix}$$

i) $x=0$ $y=0$ $\underline{\Sigma}^{\text{post}} = \begin{bmatrix} 1 & 0 \\ 0 & 1/2 \end{bmatrix}$ $\underline{\mu}^{\text{post}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

ii) $x=1$ $y=0$ $\underline{\Sigma}^{\text{post}} = \begin{bmatrix} 2/3 & -1/3 \\ -1/3 & 2/3 \end{bmatrix}$ $\underline{\mu}^{\text{post}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

iii) $x=100$ $y=100$ $\underline{\Sigma}^{\text{post}} \approx \begin{bmatrix} 2/100^2 & -1/100 \\ -1/100 & 1 \end{bmatrix}$ $\underline{\mu}^{\text{post}} \approx \begin{bmatrix} 1 \\ 1/100 \end{bmatrix}$

i) data doesn't tell us anything about m : $y = mx + c$ & $x=0$

This is why posterior mean & variance of m stays @ prior (mean 0 Variance 1)

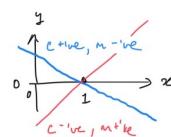
data tells us that c is likely to be ≈ 0 , but the observation noise is quite large

so only weak evidence for this. This why posterior mean is 0 and posterior variance is $1/2$ (prior variance was 1).

ii) The data can be explained by either positive c & negative m or negative c & positive m :

Hence posterior covariance between m & c is negative ($-1/3$)

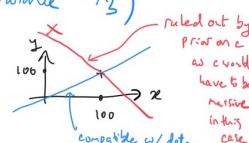
Dataset provides weak evidence that m & c are both small in magnitude & due to symmetry both still have posterior mean 0 (uncertainty has reduced from variance 1 to variance $2/3$)



iii) The data $y=100$ $x=100$ rules out negative m :

We are therefore quite certain that m is close to 1 (posterior mean very close to this value & variance $\approx 2/100^2$)

The data give almost no information about c though with posterior mean & variance close to the prior -



You might like to compute the predictive mean for these three cases i.e. the mean of $p(y^*|x^*, x, y)$ where x^* is the location of a new test input and y^* is the corresponding output.

Predictive mean : use $y^* = Mx^* + C + \varepsilon^*$ now average wrt $p(M, C, \varepsilon^* | x, y)$:

$$\mathbb{E}_{\substack{p(M, C, \varepsilon^* | x, y) \\ \text{independent from } x^* \text{ & } y}} [Mx^* + C + \varepsilon^*] = \mu_M^{\text{post}} x^* + \mu_C^{\text{post}} + b$$

i) predictive mean = 0 (horizontal line)

ii) predictive mean = 0 (horizontal line)

iii) predictive mean = $\frac{1}{100^2 + 2} \left(100^2 x^* + 100 \right)$ (line with gradient v. close to 1 and a small positive intercept)

Classification

2. Probit Classification

Consider classification as described in lectures, but with the following model

$$y^{(n)} = H(\mathbf{w}^\top \mathbf{x}^{(n)} + \epsilon_n)$$

where ϵ_n is Gaussian with mean 0 and variance σ^2 and $H(\cdot)$ is the Heaviside step function.

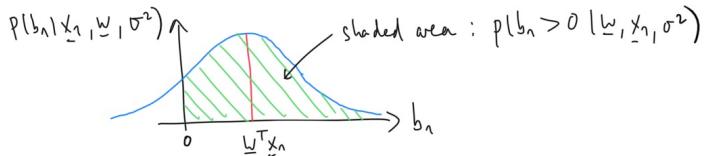
- (a) Compute the probability $P(y^{(n)} = 1 | \mathbf{x}_n, \mathbf{w}, \sigma^2)$ in terms of the Gaussian cumulative distribution. Sketch $P(y^{(n)} = 1 | \mathbf{x}_n, \mathbf{w}, \sigma^2)$ as a function of the inputs \mathbf{x}_n in the case where they are one dimensional.

- (b) What happens as the noise variance tends to infinity $\sigma^2 \rightarrow \infty$?

a) $y_n = H(\underline{\mathbf{w}}^\top \underline{\mathbf{x}}_n + \varepsilon_n) \quad \varepsilon_n \sim N(0, \sigma^2)$

let $b_n = \underline{\mathbf{w}}^\top \underline{\mathbf{x}}_n + \varepsilon_n \quad b_n \sim N(\underline{\mathbf{w}}^\top \underline{\mathbf{x}}_n, \sigma^2)$

$\therefore P(y_n=1 | \underline{\mathbf{w}}, \underline{\mathbf{x}}_n, \sigma^2) = P(b_n > 0 | \underline{\mathbf{w}}, \underline{\mathbf{x}}_n, \sigma^2)$



$\therefore P(y_n=1 | \underline{\mathbf{w}}, \underline{\mathbf{x}}_n, \sigma^2) = \int_0^\infty N(b_n; \underline{\mathbf{w}}^\top \underline{\mathbf{x}}_n, \sigma^2) db_n$

let $v_n = -\left(\frac{b_n - \underline{\mathbf{w}}^\top \underline{\mathbf{x}}_n}{\sigma}\right)$ & transform to integral over this quantity (v_n will be distributed according to a standard normal)

$$P(y_n=1 | \underline{\mathbf{w}}, \underline{\mathbf{x}}_n, \sigma^2) = \int_{-\infty}^{\frac{\underline{\mathbf{w}}^\top \underline{\mathbf{x}}_n}{\sigma}} N(v_n; 0, 1) dv_n$$

$$P(y_n=1 | \underline{\mathbf{w}}, \underline{\mathbf{x}}_n, \sigma^2) = \text{CDF}\left(\frac{\underline{\mathbf{w}}^\top \underline{\mathbf{x}}_n}{\sigma}\right)$$

b) $\sigma^2 \rightarrow \infty \Rightarrow P(y_n=1 | \underline{\mathbf{x}}_n, \underline{\mathbf{w}}, \sigma^2) = 1/2$

The noise swamps $\underline{\mathbf{w}}^\top \underline{\mathbf{x}}_n$ resulting in an output which is a coin toss

3. Multi-class Classification

Consider a multi-class classification problem with K classes. The training labels are represented by K dimensional vectors \mathbf{t}_n which has a single element set to 1, indicating the class membership, and all other values are set to 0. The inputs are multi-dimensional vectors \mathbf{x}_n . The goal is to use a training set of input vectors and output labels $\{\mathbf{x}_n, \mathbf{t}_n\}_{n=1}^N$ to enable prediction at unseen input locations.

A friend suggests using a soft-max function for this purpose which is parameterised by weights $\mathbf{W} = \{\mathbf{w}_k\}_{k=1}^K$. The output of the function is a vector, \mathbf{y} , with elements given by

$$y_i(\mathbf{x}; \mathbf{W}) = \frac{\exp(\mathbf{w}_i^\top \mathbf{x})}{\sum_{k=1}^K \exp(\mathbf{w}_k^\top \mathbf{x})}.$$

- (a) What happens to the softmax function as the magnitude of the weights tends to infinity?
- (b) Interpretting the output of the softmax as $y_i = p(t_i = 1 | \mathbf{W}, \mathbf{x})$ write down a cost-function for training this network based on the log-probability of the training data given the weights \mathbf{W} and inputs $\{\mathbf{x}_n\}_{n=1}^N$.
- (c) What is the relationship between this network and logistic regression?

$$y_i(x; \underline{w}) = \frac{e^{\underline{w}_i^\top \underline{x}}}{\sum_{k=1}^K e^{\underline{w}_k^\top \underline{x}}}$$

let $\hat{\underline{w}}_{\max} = \arg \max_{\hat{\underline{w}}_i \in \hat{\underline{w}}_1, \dots, \hat{\underline{w}}_K} \hat{\underline{w}}_i^\top \underline{x}$

a) let $\underline{w}_i = \beta \hat{\underline{w}}_i$ unit vector in direction of \underline{w}_i
magnitude

$$y_i(x; \underline{w}) = \frac{e^{\beta \hat{\underline{w}}_i^\top \underline{x}}}{\sum_{k=1}^K e^{\beta \hat{\underline{w}}_k^\top \underline{x}}} = \frac{e^{\beta (\hat{\underline{w}}_i - \hat{\underline{w}}_{\max})^\top \underline{x}}}{\sum_{k=1}^K e^{\beta (\hat{\underline{w}}_k - \hat{\underline{w}}_{\max})^\top \underline{x}}}$$

$(\hat{\underline{w}}_i - \hat{\underline{w}}_{\max})^\top \underline{x}$ are a set of K scalars, one of which is zero & the others are negative

$$\therefore \text{as } \beta \rightarrow \infty \quad e^{\beta (\hat{\underline{w}}_i - \hat{\underline{w}}_{\max})^\top \underline{x}} \rightarrow \begin{cases} 0 & \text{if } \hat{\underline{w}}_i \neq \hat{\underline{w}}_{\max} \\ 1 & \text{if } \hat{\underline{w}}_i = \hat{\underline{w}}_{\max} \end{cases}$$

$$\Rightarrow y_i(x; \underline{w}) \rightarrow \begin{cases} 0 & \text{if } i \neq i_{\max} \\ 1 & \text{if } i = i_{\max} \end{cases}$$

hence the general name "softmax"

This is a zero/one encoding of the arg max function:

$$\text{ie as } \beta \rightarrow \infty \quad y_i(x; \underline{w}) \rightarrow \prod_{k=1}^K \left(\underbrace{\left(\arg \max_k \hat{\underline{w}}_k^\top \underline{x} = i \right)}_{\text{indicator function that takes value 1 if the argument is true & 0 otherwise.}} \right)$$

NB
many people will be able to intuit this result without going through all of this.

$$b) \quad y_i = p(t_i=1 | w, x)$$

$$\begin{aligned} P\left(\left\{\underline{t}^{(n)}\right\}_{n=1}^N \mid \underline{w}, \underline{x}^{(n)}\right) &= \prod_{n=1}^N P(t^{(n)} \mid \underline{w}, \underline{x}^{(n)}) \\ &= \prod_{n=1}^N \prod_{k=1}^K y_k(x^{(n)}, w)^{t_k^{(n)}} \\ &= e^{\sum_n \sum_k t_k^{(n)} \log y_k(x^{(n)}, w)} \end{aligned}$$

\Rightarrow cost function could be

$$\underset{\underline{w}}{\text{arg max}} \sum_{n=1}^N \sum_{k=1}^K t_k^{(n)} \log y_k(x^{(n)}, w)$$

c) Consider $K=2$ (two classes)

$$\begin{aligned} y_1(x, w) &= \frac{e^{w^\top x}}{e^{w_1^\top x} + e^{w_2^\top x}} = \frac{1}{1 + e^{(w_2 - w_1)^\top x}} \\ &= \frac{1}{1 + e^{v^\top x}} \Rightarrow \text{equivalent to logistic} \\ &\quad \text{classification when there are 2 classes} \end{aligned}$$

Dimensionality Reduction

4. Principal Component Analysis[†]

A dataset comprises of N data points $\{x_n\}_{n=1}^N$. The data are zero-mean, $\mu = \frac{1}{N} \sum_{n=1}^N x_n = 0$, and the data-covariance is given by, $\Sigma = \frac{1}{N} \sum_{n=1}^N x_n x_n^\top$.

- (a) Describe the purpose of the principal component analysis algorithm and outline one procedure for computing the first principal component of a dataset.

- (b) Consider a data-covariance $\Sigma = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$. Show that the two eigenvectors for this matrix are $\frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ and $\frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}$ and state the corresponding eigenvalues. What is the first principal component when $-1 < \rho < 0$?

a) PCA carries out dimensionality reduction

Find eigenvector of covariance matrix Σ with the highest eigenvalue

$$\text{ie find } \underline{\Sigma} \underline{e}_\mu = \lambda \underline{e}_\mu \quad \text{principal component} = \underline{e}_{\mu_{\max}} \quad \lambda_{\max} = \sigma_{\max} \lambda_\mu$$

$$b) \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1+\rho \\ \rho+1 \end{bmatrix} = (1+\rho) \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\lambda = 1+\rho$$

$$\frac{1}{\sqrt{2}} \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1-\rho \\ \rho-1 \end{bmatrix} = (1-\rho) \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$\lambda = 1-\rho$$

$$\text{if } -1 < \rho < 0 \quad \text{then} \quad \lambda_{\max} = 1-\rho \quad \Rightarrow \quad \underline{e}_{\max} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

5. Principal Component Analysis and Auto-Encoders**

A data-scientist would like to summarise high dimensional data points \mathbf{y}_n in terms of a single scalar variable x_n . They use an encoding weight \mathbf{w} to produce the summary, $x_n = \mathbf{w}^\top \mathbf{y}_n$, and a decoding weight \mathbf{r} to reconstruct the data point from the summary, $\hat{\mathbf{y}}_n = \mathbf{r}x_n$. This architecture is called an auto-encoder. The data-scientist would like to learn the encoding and decoding weights by optimising the squared error of the reconstruction,

$$\mathcal{C} = \sum_n \|\mathbf{y}_n - \hat{\mathbf{y}}_n\|^2.$$

- (a) Minimise the cost \mathcal{C} with respect to the decoding weights \mathbf{r} , returning an expression for them in terms of x_n and \mathbf{y}_n .
- (b) Substitute your expression for the optimised decoding weights \mathbf{r} into \mathcal{C} to obtain the cost purely in terms of the encoding weights \mathbf{w} .
- (c) Now consider minimising the cost derived in part (b) with respect to the encoding weights. What is the solution? Is it unique?

It may be useful to know that the solution to the optimisation problem $\mathbf{z}^* = \arg \max_{\mathbf{z}} \frac{\mathbf{z}^\top H \mathbf{z}}{\mathbf{z}^\top \mathbf{z}}$ is the largest eigenvector of matrix H (arbitrarily scaled), $\mathbf{z}^* \propto \mathbf{e}_1$.

$$\begin{aligned}
 & \boxed{\quad} \quad \mathcal{C}(\mathbf{w}, \Sigma) = \sum_n \sum_d (y_{dn} - r_d x_n)^2 \\
 & \frac{\partial \mathcal{C}(\mathbf{w}, \Sigma)}{\partial r_d} = \sum_n 2(y_{dn} - r_d x_n) x_n \\
 & \Rightarrow r_d = \frac{\sum_n y_{dn} x_n}{\sum_n x_n^2} \quad (\text{just like linear regression})
 \end{aligned}$$

$$\begin{aligned}
 C(\underline{w}, \underline{\Gamma}(\underline{w})) &= \sum_n \sum_d \left(y_{nd} - \frac{\sum_n y_{nd} x_n \cdot x_n}{\sum_a x_a^2} \right)^2 \\
 &= \sum_{nd} y_{nd}^2 - 2 \sum_n \sum_d \underbrace{y_{nd} y_{nd}}_{\sum_a x_a^2} + \underbrace{\sum_n \sum_d y_{nd} y_{nd} x_n \cdot x_n}_{\left(\sum_a x_a^2\right)^2} \\
 &\quad \downarrow \text{does not depend on } \underline{w} \\
 &= \sum_{nd} y_{nd}^2 - \sum_{n,d} y_{nd} y_{nd} \frac{x_n \cdot x_n}{\sum_a x_a^2} \\
 &\quad \downarrow \frac{\sum_n y_n^2}{\underline{w}^T \underline{\Sigma}_y \underline{\Sigma}_y^T \underline{w}} \\
 \Rightarrow \arg \max_{\underline{w}} \sum_{n,d} y_{nd} y_{nd} \frac{x_n \cdot x_n}{\sum_a x_a^2} &= \frac{\sum_n y_n^2 y_n \underline{w}^T \underline{y}_n \underline{w}^T \underline{y}_n}{\sum_a \underline{w}^T \underline{y}_n \underline{y}_n \underline{w}^T \underline{w}} = \frac{\underline{w}^T \underline{\Sigma}_y^2 \underline{w}}{\underline{w}^T \underline{\Sigma}_y \underline{w}}
 \end{aligned}$$

c) Now we $\underline{v} = \underline{\Sigma}_y^{1/2} \underline{w}$ & solve for \underline{v} to find \underline{w} via $\underline{w} = (\underline{\Sigma}_y^{1/2})^{-1} \underline{v}$

$$\begin{aligned}
 \arg \max_{\underline{v}} \frac{\underline{v}^T \underline{\Sigma}_y \underline{v}}{\underline{v}^T \underline{v}} &\Rightarrow \underline{v} \propto \text{largest eigenvector of } \underline{\Sigma}_y = \underline{e}_1 \\
 &\Rightarrow \underline{w} \propto \text{largest eigenvector of } \underline{\Sigma}_y \text{ too} = \underline{e}_1 \\
 \text{as } \underline{\Sigma}_y &= \underline{E} \underline{\Lambda} \underline{E}^T \xrightarrow{\text{matrix of eigenvectors}} \text{share eigenvectors} \\
 \underline{\Sigma}_y^{1/2} &= \underline{E}^T \underline{\Lambda}^{1/2} \underline{E} \xrightarrow{\text{share eigenvectors}} \\
 \underline{\Sigma}_y^{-1/2} &= \underline{E}^T \underline{\Lambda}^{-1/2} \underline{E} \xrightarrow{\Rightarrow \underline{w} \propto \underline{\Sigma}_y^{-1/2} \underline{e}_1 \propto \underline{e}_1}
 \end{aligned}$$

Not unique or free to scale \underline{w} arbitrarily
& rescale $\underline{\Gamma}$ accordingly to compensate

Clustering

6. K-means clustering*

Consider the K-means algorithm that seeks to minimise the cost function

$$C = \sum_{n=1}^N \sum_{k=1}^K s_{nk} \|x_n - m_k\|^2$$

where m_k is the mean (centre) of cluster k , x_n is the data point n , $s_{nk} = 1$ indicates that the n th data point has been assigned to cluster k , and $s_{nk} = 0$ indicates that it has not been assigned to that cluster. There are N data points and K clusters.

- (a) Given all the assignments $\{s_{nk}\}$ determine the value of m_k that minimises the cost C and give an interpretation in terms of the K-means algorithm.
- (b) You would like to automatically learn the number of clusters K from data. One possibility is to minimise the cost C as a function of K . Explain whether this is a good idea or not, and what the solution to this minimisation is.
- (c) In many real-world applications, data points arrive sequentially and one wants to cluster them as they come in. Devise a sequential variant of the k-means algorithm which takes in one data point at a time and updates the means $\{m_1, \dots, m_K\}$ sequentially without revisiting previous data points. Describe your sequential algorithm.

$$C(\{S\}, \{M\}) = \sum_{n=1}^N \sum_{k=1}^K S_{nk} \|x_n - m_k\|^2 = \sum_n \sum_{k=1}^K S_{nk} (x_{nd} - m_{kd})^2$$

$$\frac{\partial C(\{S\}, \{M\})}{\partial M_{le}} = -2 \sum_{n=1}^N S_{nl} (x_{ne} - m_{le}) = 0$$

$$\Rightarrow \sum_{n=1}^N S_{nl} x_{ne} = M_{le} \sum_{n=1}^N S_{nl}$$

$$\Rightarrow M_{le} = \frac{\sum_{n=1}^N S_{nl} x_{ne}}{\sum_{n=1}^N S_{nl}}$$

sum of all of the data values assigned to k^{th} cluster

Or in vector form $\underline{m}_k = \frac{\sum_{n=1}^N S_{nk} \underline{x}_n}{\sum_{n=1}^N S_{nk}}$

number of datapoints assigned to k^{th} cluster

b) Bad idea : if $K=N$ & $\underline{m}_n = \underline{x}_n$ cost is zero
 i.e. each datapoint is a "cluster"



which is lowest possible

c) Several possible options for this are but a logical soln would be:

$$\text{Since } \underline{M}_k = \frac{\sum_{n=1}^N \underline{x}_n s_{nk}}{\sum_{n=1}^N s_{nk}} = \frac{\sum_{n=1}^{N-1} \underline{x}_n s_{nk} + \underline{x}_N s_{nk}}{\sum_{n=1}^{N-1} s_{nk} + s_{nk}}$$

$\beta_k^{(N-1)}$

$\alpha_k^{(N-1)}$

New datum arrives

- assign to nearest cluster and set assignment vector $\{\hat{s}_{nk}\}_{k=1}^K$

- update statistics $\beta_k^{(N)} = \beta_k^{(N-1)} + \underline{x}_N s_{nk}$ [running sum of data vectors associated w/ kth cluster]

$\alpha_k^{(N)} = \alpha_k^{(N-1)} + s_{nk}$ [running count of # of data points associated w/ kth cluster]

- recompute mean $\underline{M}_k = \frac{\beta_k^{(N)}}{\alpha_k^{(N)}}$

- repeat

7. The KL Divergence

The KL Divergence between two discrete distributions $p(x = k) = p_k$ and $p(x = k) = q_k$ is defined as $\mathcal{KL}(q, p) = \sum_{k=1}^K q_k \log \frac{q_k}{p_k}$.

- (a) Prove that the KL Divergence is non-negative and that it attains its unique minimum when $q_k = p_k$.
- (b) A machine learner has a target distribution $\mathbf{p} = [p_1, p_2, p_3, p_4, p_5, p_6] = [1, 1, 0, 0, 1, 1]/4$ which they want to fit with approximating distribution q . The choices for q available to the machine learner are: $q_1 = [1, 1, 0, 0, 0, 0]/2$, $q_2 = [0, 1, 1, 0, 0, 0]/2$, $q_3 = [0, 0, 1, 1, 0, 0]/2$, $q_4 = [0, 0, 0, 1, 1, 0]/2$, $q_5 = [0, 0, 0, 0, 1, 1]/2$, and $q_6 = [1, 1, 1, 1, 1, 1]/6$.
 - i. Determine the distribution(s) q_i that minimises $\mathcal{KL}(q_i, p)$. Comment on your result.
 - ii. Determine the distribution(s) q_i that minimises $\mathcal{KL}(p, q_i)$. Comment on your result.

Note that, by convention, $0 \times \log 0 = 0$ since $\lim_{\Delta \rightarrow 0} \Delta \log \Delta = 0$.

$$a) \quad \mathcal{KL}(q || p) = \sum_k q_k \log \frac{q_k}{p_k}$$

constraint that probabilities sum to 1

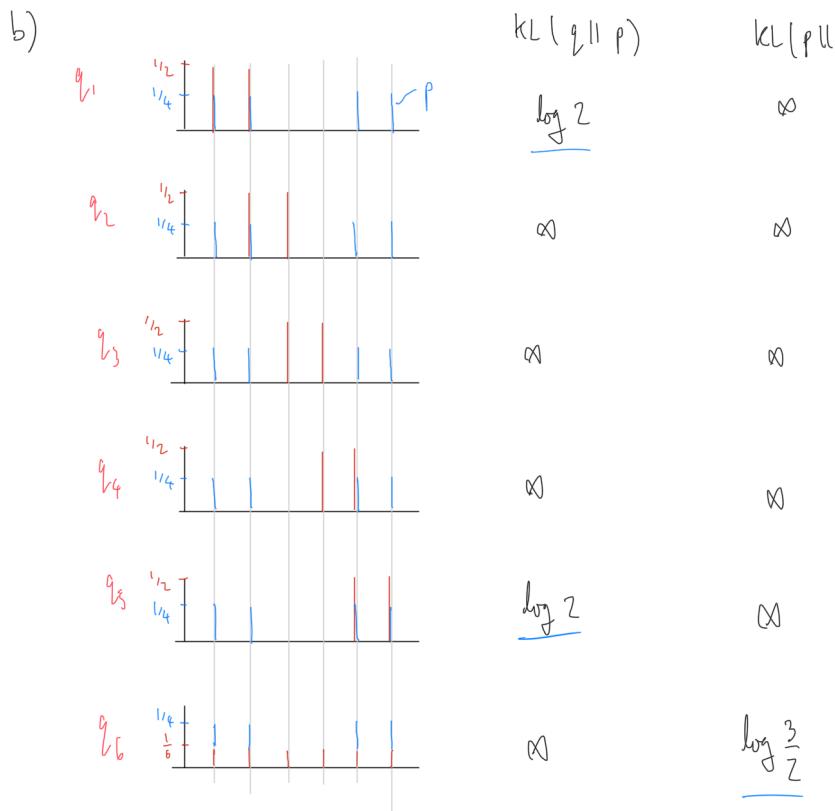
$$\mathcal{O} = \frac{\partial}{\partial q_l} \left[\sum_k q_k \log \frac{q_k}{p_k} + \lambda \left(\sum_k q_k - 1 \right) \right]$$

↑
Lagrange multiplier

$$\mathcal{O} = \log \frac{q_l}{p_l} + \frac{q_l}{q_l} + \lambda$$

$$\Rightarrow q_l = p_l \text{ @ which } \mathcal{KL}(q^{\text{opt}} || p) = \sum_k p_k \log \frac{p_k}{p_k} = 0$$

$$\frac{\partial^2}{\partial q_l^2} \mathcal{KL}(q || p) = \frac{1}{q_l} > 0 \Rightarrow \text{maximum}$$



$KL(q||p)$ is zero avoiding & tends to fit single modes of p with q

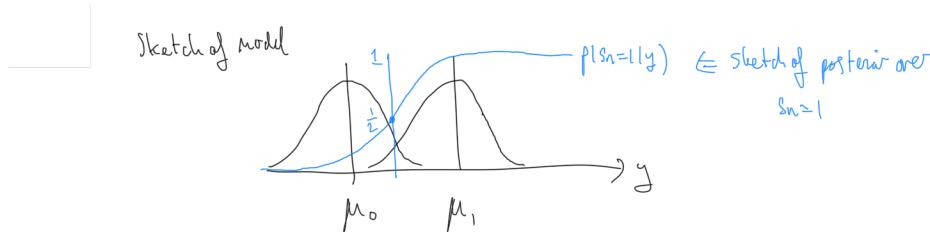
$KL(p||q)$ is not zero avoiding & tends to cover multiple modes of p with q

8. Mixtures of Gaussians and EM*

A set of N scalar data points $\{y_n\}_{n=1}^N$ are modelled using a mixture of Gaussians containing two equiprobable components with unknown means (μ_0 and μ_1) and unit variances,

$$p(s_n = 1) = \frac{1}{2}, \quad p(y_n | s_n = 0) = \mathcal{N}(y_n; \mu_0, 1), \quad p(y_n | s_n = 1) = \mathcal{N}(y_n; \mu_1, 1).$$

- (a) Compute the posterior distribution over the components, $p(s_n = 1 | y_n)$ and sketch how this varies as a function of the observed data y_n . Briefly discuss how this relates to logistic classification.
- (b) Explain how your solution to (a) can be used in the EM algorithm to estimate the component means. Your answer should include a expression for the M-step update.
- (c) Do you expect the EM algorithm to overfit when used to train this model?



$$\begin{aligned}
 a) \quad p(s_n = 1 | y_n) &= \frac{p(s_n = 1, y_n)}{p(s_n = 1, y_n) + p(s_n = 0, y_n)} = \frac{1}{1 + \frac{p(s_n = 0, y_n)}{p(s_n = 1, y_n)}} \\
 &= \frac{1}{1 + \exp \left(-\frac{1}{2} (y_n - \mu_0)^2 + \frac{1}{2} (y_n - \mu_1)^2 \right)} \\
 &= \frac{1}{1 + \exp \left(-\underbrace{[y_n(\mu_1 - \mu_0) + \frac{\mu_1^2 - \mu_0^2}{2}]}_{\alpha} \right)} \\
 &= \frac{1}{1 + \exp(-\alpha)} = \text{logistic function}
 \end{aligned}$$

Just like logistic regression with weights $(\mu_1 - \mu_0)$ & bias $\frac{\mu_1^2 - \mu_0^2}{2}$

b) E-Step computes $p(s_n=1|y_n)$ using expression above

$$\underline{\text{M-Step}} \quad \mu_1^{(\text{new})} = \frac{\sum_n p(s_n=1|y_n) y_n}{\sum_n p(s_n=1|y_n)}$$

mean of data
 each weighted
 by probability
 they came
 from class 1

$$\mu_0^{(\text{new})} = \frac{\sum_n (1 - p(s_n=1|y_n)) y_n}{\sum_n (1 - p(s_n=1|y_n))}$$

p($s_n=0|y_n$)
 as above, but
 weighted by
 prob came from
 class 0

c) Overfitting in a MoG model can occur when variances shrink to zero. As the variances are fixed to unity in this model over fitting is less likely.

The EM Algorithm

9. Factor Analysis and EM*

The noisy depth sensor from question 2 in example sheet 1 is used to collect measurements of the distances to a set of N objects that are unknown distances d_n metres away. The object depths can be assumed, *a priori*, to be distributed according to independent standard Gaussian distributions $p(d_n) = \mathcal{N}(d_n; 0, 1)$. As before, the depth sensor returns y_n a noisy measurement of the depth, that is also assumed to be Gaussian $p(y_n|d_n, \sigma_y^2) = \mathcal{N}(y_n; d_n, \sigma_y^2)$.

The variance of the σ_y^2 sensor noise is unknown and must be estimated from the measured depths $\{y_n\}_{n=1}^N$ using maximum likelihood.

- (a) Derive the steps of the EM algorithm for performing this. Your answer should include the E-Step (you may find your solution to question 2 in example sheet 1 useful here) and the M-Step.
- (b) An alternative approach to maximum-likelihood learning would optimise the log-likelihood $p(\{y_n\}_{n=1}^N | \sigma_y^2)$ directly. Compute the objective function for this procedure. How do you think this approach will perform compared to EM?

$$p(d_n) = N(d_n; 0, 1)$$

$$p(y_n | d_n, \sigma_y^2) = N(y_n; d_n, \sigma_y^2)$$

in E-Step set $q_k(d_n) = p(d_n | y_n, \sigma_y^2) = N(d_n; \mu_{d_n|y_n}, \sigma_{d_n|y_n}^2)$

where $\mu_{d_n|y_n} = \frac{y_n}{1 + \sigma_y^2}$ $\sigma_{d_n|y_n}^2 = \frac{\sigma_y^2}{1 + \sigma_y^2}$ (see Q2 Examples Sheet 1)

in M-Step

$$\sigma_y^2 = \underset{\sigma_y^2}{\operatorname{arg\ max}} \underbrace{\mathbb{E}_q \left[\log p(\{d_n\}_{n=1}^N, \{y_n\}_{n=1}^N | \sigma_y^2) \right]}_{G(\sigma_y^2)}$$

↑ "average log-joint"

$$G(\sigma_y^2) = \mathbb{E}_q \left[\sum_n \log p(d_n) + \sum_n \log p(y_n | d_n, \sigma_y^2) \right]$$

$$= C - \frac{1}{2\sigma_y^2} \sum_n \mathbb{E}_{q(d_n)} \left[(y_n - d_n)^2 \right] - \frac{N}{2} \log 2\pi \sigma_y^2 \quad \star$$

$$= C - \frac{1}{2\sigma_y^2} \sum_n (y_n^2 + \mathbb{E}_{q(d_n)} [d_n^2] - 2y_n \mathbb{E}_{q(d_n)} [d_n]) - \frac{N}{2} \log 2\pi \sigma_y^2$$

$$\frac{dG(\sigma_y^2)}{d\sigma_y^2} = \frac{1}{2\sigma_y^4} \sum_n (y_n^2 + \mathbb{E}_{q(d_n)} [d_n^2] - 2y_n \mathbb{E}_{q(d_n)} [d_n]) - \frac{N}{2\sigma_y^2} = 0$$

$$\therefore \sigma_y^2 = \frac{1}{N} \sum_n y_n^2 + \frac{1}{N} \sum_n (\text{Model}_n^2 + \sigma_{\text{noise}}^2) - \frac{2}{N} \sum_n y_n / \text{Model}_n$$

Algorithmically this is what you need, although for intuition ~~it~~ yields:

$$\sigma_y^2 = \frac{1}{N} \sum_n \mathbb{E}_{q(d_n)} [(y_n - d_n)^2]$$

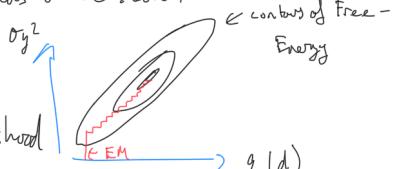
i.e. the noise variance is the average (over data points n & the posterior) of the squared error between y_n & d_n .

b) $p(y_n | \sigma_y^2) = N(y_n; 0, 1 + \sigma_y^2)$

$$\therefore \log p(\{y_n\}_{n=1}^N | \sigma_y^2) = -\frac{N}{2} \log[2\pi(1 + \sigma_y^2)] - \frac{1}{2(1 + \sigma_y^2)} \sum_n y_n^2$$

Optimising this function for σ_y^2 (eg by gradient ascent) would be much faster than EM. EM performs co-ordinate ascent of the free-energy in q & θ which tends to be slow.

It turns out that if EM is analytic, so too is direct optimisation of the log likelihood which is almost always to be preferred.



10. Gaussian Mixture Models and EM** (beyond tripos standard; optional)

A Gaussian Mixture Model for D -dimensional data $\{\mathbf{x}_n\}_{n=1}^N$ comprises a categorical distribution over the class membership variables $p(s_n = k|\theta) = \pi_k$ and a general multivariate Gaussian distribution over the observed data given the class membership variables, $p(\mathbf{x}_n|s_n = k, \theta) = \mathcal{N}(\mathbf{x}_n; \mathbf{m}_k, \Sigma_k)$, (i.e. Σ_k is not isotropic). The posterior distribution over the class membership variables has been computed $p(s_n = k|\mathbf{x}_n, \theta) = q_{n,k}$.

- (a) Derive the M-Step update of the EM algorithm, i.e. the formulae for updating the parameters $\{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$ using the posterior probabilities, $q_{n,k}$.
- (b) A friend suggests that it might be possible to speed up the EM algorithm by updating the posterior distribution of only a subset of $K < N$ data points during the E-Step that are selected uniformly at random each time, and then updating the parameters using the same expressions derived in part (a).
 - i. Will this procedure converge to a (local) optimum of the likelihood?
 - ii. Will this partial E-Step update procedure be computationally more efficient?

You may find the following identities useful,

$$\frac{d}{d\alpha} \log \det(\Sigma) = \text{trace} \left(\Sigma^{-1} \frac{d\Sigma}{d\alpha} \right), \quad \frac{d}{d\alpha} \Sigma^{-1} = -\Sigma^{-1} \frac{d\Sigma}{d\alpha} \Sigma^{-1}.$$

$$\boxed{\quad \quad \quad \mathcal{F}(\underline{q}(\underline{s}), \theta) = \sum_n q(s_n) \frac{\log p(x_n | s_n, \theta) p(s_n | \theta)}{q(s_n)}}$$

$$q(s_n = k) = q_{nk}$$

$$\therefore \mathcal{F}(\underline{q}(\underline{s}), \theta) = \sum_n \sum_k q_{nk} \frac{\log p(z_n | s_n = k, \theta) p(s_n = k | \theta)}{q(s_n = k)}$$

where

$$p(x_n | s_n = k, \theta) = N(x_n; \mu_k, \Sigma_k) \quad p(s_n = k | \theta) = \pi_k$$

$$\therefore \mathcal{F}(\underline{q}(\underline{s}), \theta) = \sum_n \sum_k q_{nk} \left[-\frac{1}{2} \log \det(\Sigma_k) - \frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) + \log \pi_k \right] + C$$

$$\begin{aligned} \frac{\partial \mathcal{F}}{\partial \mu_{k'd'}} &= \frac{1}{2} \sum_n \sum_k q_{nk} \sum_{ij} (x_{in} - \mu_{ki}) \Sigma_{kij}^{-1} (x_{jn} - \mu_{kj}) = 0 \\ &= \frac{1}{2} \sum_n q_{nk} \left[\sum_j \Sigma_{k'd'j}^{-1} (x_{jn} - \mu_{kj}) + \sum_i (x_{in} - \mu_{ki}) \Sigma_{k'd'}^{-1} \right] \\ &= \sum_n q_{nk} \sum_j \Sigma_{k'd'j}^{-1} (x_{jn} - \mu_{kj}) = \left[\sum_{k'}^{-1} \sum_n q_{nk} (x_n - \mu_{k'}) \right]_{d'} \end{aligned}$$

does not depend on parameters

$$\therefore \underline{\mu}_n = \frac{\sum_n q_{nk} \underline{x}_n}{\sum_n q_{nk}}$$

alternatively, you could differentiate w.r.t Σ_{kij}^{-1} & use $\det(\Sigma) = \frac{1}{\det(\Sigma^{-1})}$

$$\frac{d\varphi}{d\Sigma_{kij}} = \frac{1}{\det(\Sigma)} \left[\sum_{n,k} q_{nk} \left[-\frac{1}{2} \log \det(\Sigma) - \frac{1}{2} (\underline{x}_n - \underline{\mu}_k)^T \Sigma^{-1} (\underline{x}_n - \underline{\mu}_k) \right] \right]$$

$$\frac{d \log(\det \Sigma)}{d\alpha} = \text{trace} \left(\Sigma^{-1} \frac{d\Sigma}{d\alpha} \right) \quad \frac{d \Sigma^{-1}}{d\alpha} = - \Sigma^{-1} \frac{d\Sigma}{d\alpha} \Sigma^{-1}$$

$$\frac{d\varphi}{d\Sigma_{kij}} = \sum_{n,k} q_{nk} \left[-\frac{1}{2} \sum_{ij} \Sigma_{kij}^{-1} \frac{d\Sigma_{kij}}{d\Sigma_{k'j'}} + \frac{1}{2} \sum_{ijlm} (\underline{x}_{ni} - \underline{\mu}_{ki}) \Sigma_{k'l}^{-1} \frac{d\Sigma_{k'm} \Sigma_{m'l}^{-1}}{d\Sigma_{k'j'}} (\underline{x}_{nj} - \underline{\mu}_{kj}) \right]$$

$$= \sum_{n,k} q_{nk} \left[-\frac{1}{2} \sum_{ij} \Sigma_{kij}^{-1} \delta_{k'i} \delta_{i'j} \delta_{j'i} + \frac{1}{2} \sum_{ijlm} (\underline{x}_{ni} - \underline{\mu}_{ki}) \Sigma_{k'l}^{-1} \delta_{k'm} \delta_{m'l} \delta_{l'n} \delta_{nj} \Sigma_{nj}^{-1} \right]$$

$$\Rightarrow \sum_n q_{nk} \left(-\Sigma_{k'l}^{-1} + \sum_{i'j'} (\underline{x}_{ni} - \underline{\mu}_{ki}) (\underline{x}_{nj} - \underline{\mu}_{kj})^T \Sigma_{nj}^{-1} \right) = 0$$

$$= \frac{1}{2} \sum_n q_{nk} \left(-\sum_{k'j'i'} \delta_{k'i} \delta_{i'j} \delta_{j'i} + \sum_{ij} (\underline{x}_{ni} - \underline{\mu}_{ki}) \Sigma_{k'i}^{-1} \Sigma_{k'j'i'}^{-1} (\underline{x}_{nj} - \underline{\mu}_{kj}) \right)$$

$$0 = \sum_n q_{nk} \left(-\Sigma_{k'l}^{-1} + \sum_{i'j'} (\underline{x}_{ni} - \underline{\mu}_{ki}) (\underline{x}_{nj} - \underline{\mu}_{kj})^T \Sigma_{nj}^{-1} \right)$$

$$\Rightarrow \Sigma_{k'l} = \sum_n q_{nk} (\underline{x}_{ni} - \underline{\mu}_{ki}) (\underline{x}_{nj} - \underline{\mu}_{kj})^T / \sum_n q_{nk}$$

pre & post multiply by $\Sigma_{k'l}$

To find Π we need to use Lagrange multipliers :

$$\frac{d}{d\Pi_{k'l}} \left(\varphi(q, \theta) - \lambda \left(\sum_k \Pi_{k'l} - 1 \right) \right) = \frac{d}{d\Pi_{k'l}} \left(\sum_k \left(q_{nk} \log \Pi_{k'l} - \lambda \Pi_{k'l} \right) - 1 \right)$$

$$= \sum_n q_{nk} \left(\frac{1}{\Pi_{k'l}} - \lambda \right) = 0$$

$$\Rightarrow \Pi_{k'l} = \frac{\sum_n q_{nk}}{\sum_n q_{nk}} \quad \text{from constraint that } \sum_k \Pi_{k'l} = 1$$

$$= \frac{1}{N} \sum_n q_{nk}$$

- b) i) The tree-energy is guaranteed to improve (or stay the same) if only a fraction of the datapoint posteriors are updated.
 - o. The procedure will converge to a local optimum of the likelihood as for Normal EM
- ii) It's not completely clear that the procedure will help in all cases - for although the cost of the E-Step is reduced (from $O(N)$ to $O(K)$) convergence will take more steps in general. However, when there are lots of data, you don't need to visit all of them to figure out a good update for the parameters so this means the procedure is likely to be substantially more efficient in these cases.