

Module MLMI1: Introduction to Machine Learning

Example Sheet 1: Introductory Inference Problems, Bayesian Decision Theory, Regression and Classification

*Straightforward questions are marked †*

*Hard questions are marked \**

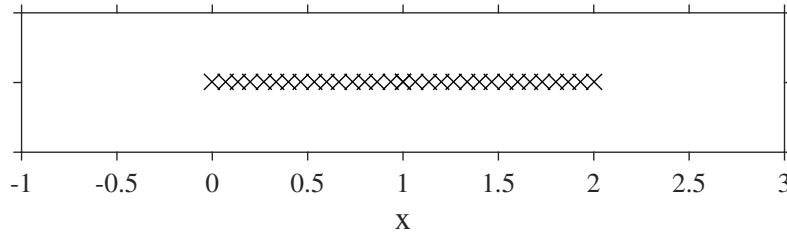
*Introductory Inference Problems*

1. Maximum likelihood fitting of a Gaussian

- (a) Explain the terms likelihood function, prior probability distribution, and posterior probability distribution, in the context of the inference of parameters  $\theta$  from data  $\mathcal{D}$ .
- (b) A random variable  $x$  is believed to have a probability distribution which is Gaussian with mean  $\mu$  and standard deviation equal to 1,

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

A sample of  $N = 32$  data points is collected  $\{x_n\}_{n=1}^N$  that are believed to be drawn independently from this distribution. The dataset is shown below:



The first and second moments of these data are  $\frac{1}{N} \sum_{n=1}^N x_n = 1$  and  $\frac{1}{N} \sum_{n=1}^N x_n^2 = 1.3$ .

Sketch the likelihood as a function of  $\mu$  for the dataset. Label the position of the maximum and its width. You do not need to compute the value of the likelihood at its maximum.



## 2. Inference in a Gaussian model

A noisy depth sensor measures the distance to an object an unknown distance  $d$  metres away. The depth can be assumed, *a priori*, to be distributed according to a standard Gaussian distribution  $p(d) = \mathcal{N}(d; 0, 1)$ . The depth sensor returns  $y$  a noisy measurement of the depth, that is also assumed to be Gaussian  $p(y|d, \sigma_y^2) = \mathcal{N}(y; d, \sigma_y^2)$ .

- || (a) Compute the posterior distribution over the depth given the observation,  $p(d|y, \sigma_y^2)$ .
- (b) What happens to the posterior distribution as the measurement noise becomes very large  $\sigma_y^2 \rightarrow \infty$ ? Comment on this result.

The formula for the probability density of a Gaussian distribution of mean  $\mu$  and variance  $\sigma^2$  is given by

$$\mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$

$$p(d) = \mathcal{N}(d; 0, 1) \quad p(y|d, \sigma_y^2) = \mathcal{N}(y; d, \sigma_y^2)$$

$$\Rightarrow p(d|y, \sigma_y^2) = \frac{p(d) p(y|d, \sigma_y^2)}{p(y)} = \mathcal{N}\left(d; \mu_{d|y}, \sigma_{d|y}^2\right)$$

LHS      RHS

*Left side is Gaussian in  $d$*

$$\text{RHS} \propto \exp\left(-\frac{1}{2\sigma_{d|y}^2}(d - \mu_{d|y})^2\right) = \exp\left(-\frac{1}{2\sigma_{d|y}^2}d^2 + \frac{d\mu_{d|y}}{\sigma_{d|y}^2} - \frac{1}{2}\frac{\mu_{d|y}^2}{\sigma_{d|y}^2}\right)$$

$$\text{LHS} \propto \exp\left(-\frac{1}{2}d^2 - \frac{1}{2\sigma_y^2}(y - d)^2\right)$$

$$= \exp\left(-\frac{1}{2}d^2\left(1 + \frac{1}{\sigma_y^2}\right) + \frac{1}{2\sigma_y^2} \cancel{y^2} - \cancel{2yd} + \frac{1}{2\sigma_y^2}\right)$$

$$\frac{1}{\sigma_{d|y}^2} = 1 + \frac{1}{\sigma_y^2} \Rightarrow \sigma_{d|y}^2 = \frac{1}{1 + 1/\sigma_y^2}$$

$$\frac{\mu_{d|y}}{\sigma_{d|y}^2} = \frac{y}{\sigma_y^2} \Rightarrow \mu_{d|y} = \sigma_{d|y}^2 \frac{y}{\sigma_y^2} = \frac{1}{1 + \sigma_y^2} y$$

$$\sigma_y^2 \rightarrow \infty \quad \sigma_{d|y}^2 \rightarrow 1 \quad \mu_{d|y} \rightarrow 0 \quad || \quad \sigma_y^2 \rightarrow 0 \quad \sigma_{d|y}^2 \rightarrow 0 \quad \mu_{d|y} \rightarrow y$$



### 3. Bayesian inference for a biased coin\*

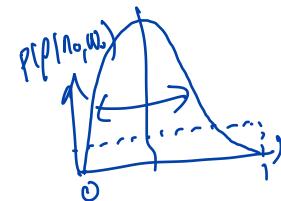
Matlab notation  $(x_1, x_2, \dots, x_N)$

A sequence of coin tosses are observed from a biased coin  $x_{1:N} = \{0, 1, 1, 0, 1, 1, 1, 1, 0\}$

where  $x_n = 1$  indicates flip  $n$  was a head and  $x_n = 0$  indicates that it was tails. An experimenter would like to estimate the coin's probability of landing heads,  $\rho$ , from these data.

The experimenter assumes that the coin flips are drawn independently from a Bernoulli distribution  $p(x_n|\rho) = \rho^{x_n}(1-\rho)^{1-x_n}$  and uses a prior distribution of the form

$$p(\rho|n_0, N_0) = \frac{1}{Z(n_0, N_0)} \rho^{n_0} (1-\rho)^{N_0-n_0}.$$



Here  $n_0$  and  $N_0$  are parameters set by the experimenter to encapsulate their prior beliefs.  $Z(n_0, N_0)$  returns the normalising constant of the distribution as a function of the parameters,  $n_0$  and  $N_0$ .

- (a) Compute the posterior distribution over the bias  $p(\rho|x_{1:N}, n_0, N_0)$ .
- (b) Compute the *maximum a posteriori* (MAP) estimate for the bias.
- (c) Provide an intuitive interpretation for the parameters of the prior distribution,  $n_0$  and  $N_0$ . For what setting of  $n_0$  and  $N_0$  does the MAP estimate become equal to the maximum-likelihood estimate?

$$\begin{aligned}
 p(\rho | x_{1:N}, n_0, N_0) &\propto p(\rho | n_0, N_0) \prod_{n=1}^N p(x_n | \rho) \\
 &= \frac{1}{Z(n_0, N_0)} \rho^{n_0} (1-\rho)^{N_0-n_0} \prod_{n=1}^N \rho^{x_n} (1-\rho)^{1-x_n} \\
 &= \frac{1}{Z(n_0, N_0)} \rho^{\sum x_n} (1-\rho)^{N_0-n_0 + N - \sum x_n} \\
 p(\rho | x_{1:N}, n_0, N_0) &= \frac{1}{Z(n', N')} \rho^{n'} (1-\rho)^{N' - n'}
 \end{aligned}$$

$p_{ML} = \frac{n}{N}$   
 $n' = n_0 + \sum x_n$   
 $N' = N_0 + N$   
 $p_{MAP} = \frac{n'}{N'}$



#### 4. Inferential game show\*

On a game show, a contestant is told the rules as follows:

There are four doors, labelled 1, 2, 3 and 4. A single prize has been hidden behind one of them. You get to select one door. Initially your chosen door will not be opened. Instead, the gameshow host will open one of the other three doors, and *he will do so in such a way as not to reveal the prize*. For example, if you first choose door 1, he will then open one of doors 2, 3 and 4, and it is guaranteed that he will choose which one to open so that the prize will not be revealed.

At this point, you will be given a fresh choice of door: you can either stick with your first choice, or you can switch to one of the other closed doors. All the doors will then be opened and you will receive whatever is behind your final choice of door.

- (a) Imagine that the contestant chooses door 1 first; then the gameshow host opens door 4, revealing nothing behind the door, as promised. Should the contestant (a) stick with door 1, or (b) switch to door 2 or 3, or (c) does it make no difference?
- (b) Use Bayes' rule to solve the problem.



## 5. Bayesian decision theory\*

A data-scientist has computed a complex posterior distribution over a variable of interest,  $x$ , given observed data  $y$ , that is  $p(x|y)$ . They would like to return a point estimate of  $x$  to their client. The client provides the data-scientist with a reward function  $R(\hat{x}, x)$  that indicates their satisfaction with a point estimate  $\hat{x}$  when the true state of the variable is  $x$ .

- Explain how to use *Bayesian Decision Theory* to determine the optimal point estimate,  $\hat{x}$ .
- Compute the optimal point estimate  $\hat{x}$  in the case when the reward function is the negative square error between the point estimate and the true value,  $R(\hat{x}, x) = -(\hat{x} - x)^2$ . Comment on your result.
- Compute the optimal point estimate  $\hat{x}$  in the case when the reward function is the negative absolute error between the point estimate and the true value,  $R(\hat{x}, x) = -|\hat{x} - x|$ . Comment on your result.

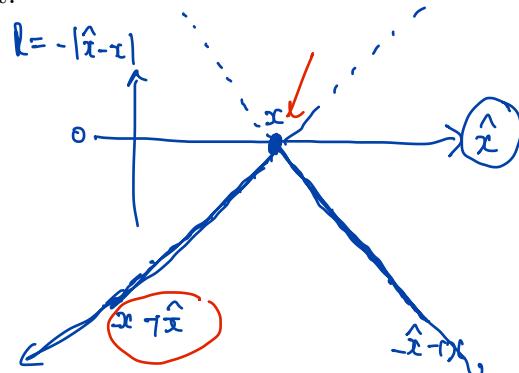
$$\hat{x}_* = \underset{\hat{x}}{\operatorname{argmax}} \int R(\hat{x}, x) p(x|y) dx$$

$$\bar{R}(\hat{x}) = \int R(\hat{x}, x) p(x|y) dx$$

$$\frac{d\bar{R}(\hat{x})}{d\hat{x}} = - \frac{d}{d\hat{x}} \int |\hat{x} - x| p(x|y) dx$$

$$+ \int (x - \hat{x})^2$$

$$\text{sign}(x - \hat{x})$$



$$\frac{dR}{d\hat{x}} = 1$$

$$\frac{dl}{dx} = -1$$

$$\text{sign}(x - \hat{x})$$



$$= \int \text{sign}(x - \hat{x}) p(x|y) dx = 0$$

$$= - \int_{-\infty}^{\hat{x}} \text{sign}(x - \hat{x}) \cdot p(x|y) dx + \int_{\hat{x}}^{\infty} +1 \cdot p(x|y) dx$$

$$\Rightarrow \hat{x} = \text{median of } p(x|y)$$

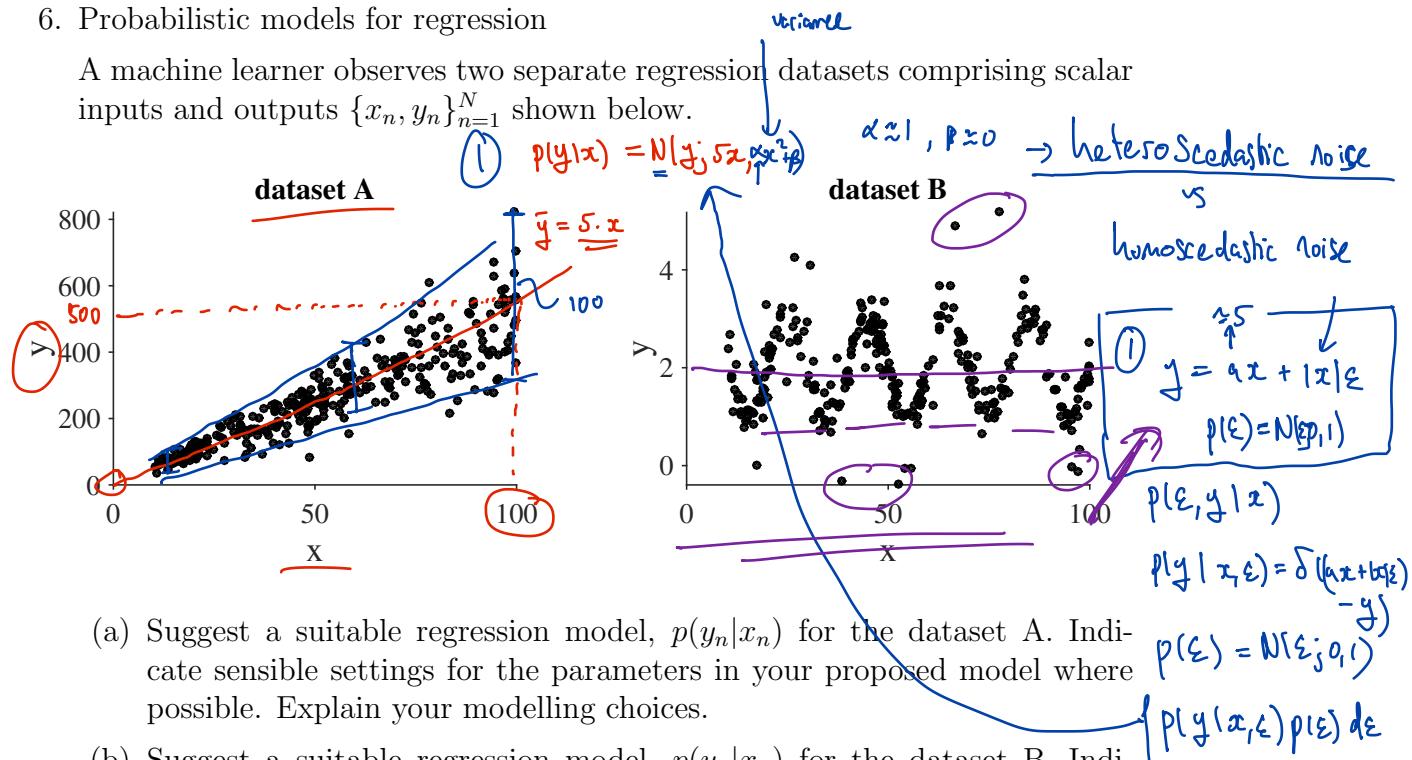


Regression

$$\begin{aligned} p(z) &= N(z; \mu, \sigma^2) \\ z &\sim N(\mu, \sigma^2) \end{aligned}$$

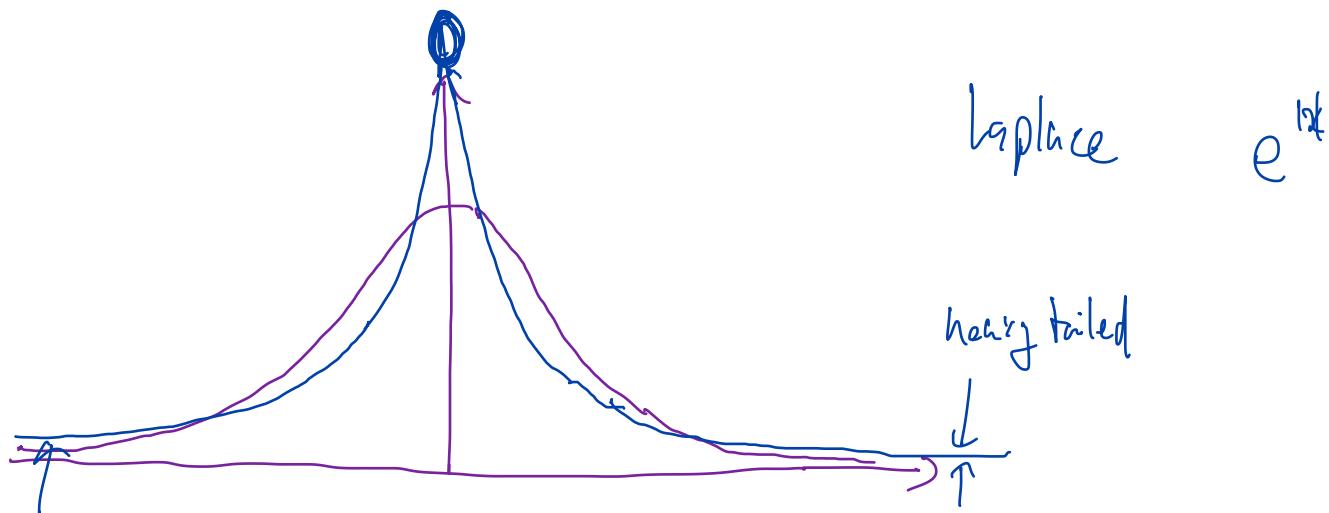
## 6. Probabilistic models for regression

A machine learner observes two separate regression datasets comprising scalar inputs and outputs  $\{x_n, y_n\}_{n=1}^N$  shown below.



- (a) Suggest a suitable regression model,  $p(y_n|x_n)$  for the dataset A. Indicate sensible settings for the parameters in your proposed model where possible. Explain your modelling choices.
- (b) Suggest a suitable regression model,  $p(y_n|x_n)$  for the dataset B. Indicate sensible settings for the parameters in your proposed model where possible. Explain your modelling choices.

$$y_n(x) = 2 + \sin\left(\frac{2\pi}{25}x_n + \phi\right) + \epsilon_n \quad \epsilon_n \sim \text{Student-t}$$



$$\begin{array}{l} \mathbb{E}(x) \\ \text{mean} \end{array}$$

$$\begin{array}{l} \mathbb{E}(x^2) \\ \text{variance} \end{array}$$

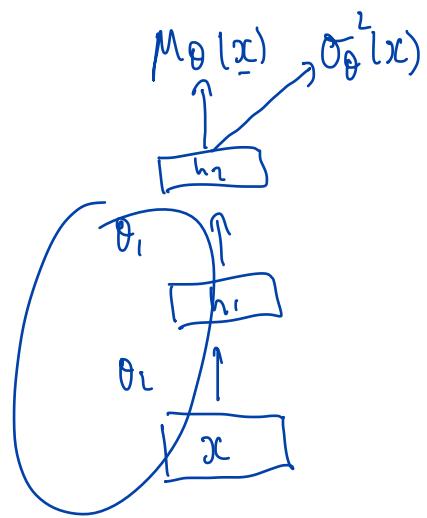
$$\begin{array}{l} \mathbb{E}(x^3) \\ \text{skew} \end{array}$$

$$\begin{array}{l} \mathbb{E}(x^4) \\ \text{variance} \end{array}$$

$-3/2 = \text{kurtosis}$  (excess)

$$f(z) = \log(1 + e^z)$$

$$p(y|x, \theta) = N(y; \mu_{\theta}(x), \sigma_{\theta}^2(x))$$





7. Maximum-likelihood learning for a simple regression model<sup>†</sup>

Consider a regression problem where the data comprise  $N$  scalar inputs and outputs,  $\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\}$ , and the goal is to predict  $y$  from  $x$ .

Assume a very simple linear model,  $y_n = ax_n + \epsilon_n$ , where the noise  $\epsilon_n$  is Gaussian with zero mean and variance 1.

- (a) Provide an expression for the log-likelihood of the parameter  $a$ .
- (b) Compute the maximum likelihood estimate for  $a$ .



## 8. Maximum-likelihood learning for multi-output regression\*

A data-scientist has collected a regression dataset comprising  $N$  scalar inputs ( $\{x_n\}_{n=1}^N$ ) and  $N$  scalar outputs ( $\{y_n\}_{n=1}^N$ ). Their goal is to predict  $y$  from  $x$  and they have assumed a very simple linear model,  $y_n = ax_n + \epsilon_n$ .

The data-scientist also has access to a second set of outputs ( $\{z_n\}_{n=1}^N$ ) that are well described by the model  $z_n = x_n + \epsilon'_n$

The noise variables  $\epsilon_n$  and  $\epsilon'_n$  are known to be zero mean correlated Gaussian variables

$$p\left(\begin{bmatrix} \epsilon_n \\ \epsilon'_n \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \epsilon_n \\ \epsilon'_n \end{bmatrix}; \mathbf{0}, \Sigma\right) \text{ where } \Sigma^{-1} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}. \Rightarrow \Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$$

- (a) Provide an expression for the log-likelihood of the parameter  $a$ .
- (b) Compute the maximum likelihood estimate for  $a$ .
- (c) Do the additional outputs  $\{z_n\}_{n=1}^N$  provide useful additional information for estimating  $a$ ? Explain your reasoning.

The formula for the probability density of a multivariate Gaussian distribution of mean  $\mu$  and covariance  $\Sigma$  is given by

$$\mathcal{N}(\mathbf{x}; \mu, \Sigma) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^\top \Sigma^{-1}(\mathbf{x} - \mu)\right).$$

$$a = \frac{\left( \sum_n x_n y_n + \underbrace{\frac{1}{2} \sum_n (z_n - \bar{x}_n) \bar{x}_n}_{\text{New bit}} \right) / \sum_n x_n^2}{\left( \sum_n x_n (\bar{y}_n + \frac{1}{2} \bar{\epsilon}'_n) \right) / \sum_n x_n^2}$$

$$\begin{aligned} \text{Correlation} &= \frac{-2/3}{4/3} = \frac{-1}{2} \\ \text{Correlation} &= \mathbb{E}(\epsilon \epsilon') \\ &= \frac{1}{\sqrt{\text{Var}(\epsilon) \text{Var}(\epsilon')}} \\ &= \frac{1}{4/3} = \frac{3}{4} \end{aligned}$$





## 9. 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  $\epsilon$  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  $p(m, c|x, y)$ .
- | (b) Show that the posterior derived in part a is consistent with the expressions for Bayesian linear regression given in lectures.
- | (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  $y = 100$

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.

$$p(m, c | y, x) = \frac{p(y|m, x, c) p(m, c | x)}{p(y|x)} \quad \left\{ \begin{array}{l} p(A|B, c) = p(B|A, c)p(A|c) \\ p(B|c) \end{array} \right.$$

$$= \frac{p(y|m, x, c) p(m)p(c)}{p(y|x)} \quad \begin{array}{l} N(y; m+x+c, 1) \\ \Downarrow \\ p(y|m, x, c) \end{array} \quad \begin{array}{l} N(m; 0, 1) \\ \Downarrow \\ p(m) \end{array} \quad \begin{array}{l} N(c; 0, 1) \\ \Downarrow \\ p(c) \end{array}$$

$B = y$   
 $A = [m]$   
 $C = x$

$$= N \left( \begin{bmatrix} m \\ c \end{bmatrix}; \underline{\mu}^{(\text{post})}, \underline{\Sigma}^{(\text{post})} \right) \quad (2)$$

$$p(m, c | y, x) \propto \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \frac{(y - (m+x+c))^2}{1}} \times \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} m^2} \times \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} c^2}$$

$$\propto e^{-\frac{1}{2} y^2 - \frac{1}{2} m^2 (x^2 + 1) - \frac{1}{2} c^2 \times 2 - mcx + y(m+x+c)}$$

$$\propto \exp \left( -\frac{1}{2} \left( \begin{bmatrix} m \\ c \end{bmatrix} - \underline{\mu}^{(\text{post})} \right)^T \left( \underline{\Sigma}^{(\text{post})} \right)^{-1} \left( \begin{bmatrix} m \\ c \end{bmatrix} - \underline{\mu}^{(\text{post})} \right) \right)$$

$$\propto \exp \left( -\frac{1}{2} \left[ \begin{bmatrix} m \\ c \end{bmatrix}^T \left( \underline{\Sigma}^{(\text{post})} \right)^{-1} \begin{bmatrix} m \\ c \end{bmatrix} + \begin{bmatrix} m \\ c \end{bmatrix}^T \left( \underline{\Sigma}^{(\text{post})} \right)^{-1} \underline{\mu}^{(\text{post})} \right] + \dots \right)$$

$$(\underline{\theta} - \underline{\mu})^T \underline{A} (\underline{\theta} - \underline{\mu}) = (\underline{\theta} - \underline{\mu})^T \underline{A} \underline{\theta} - (\underline{\theta} - \underline{\mu})^T \underline{\mu}$$

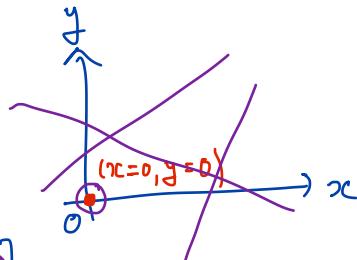
$$= \underline{\theta}^T \underline{A} \underline{\theta} - \underline{\theta}^T \underline{\mu} - \underline{\mu}^T \underline{A} \underline{\theta} + \underline{\mu}^T \underline{A} \underline{\mu}$$

$$\left( \underline{\Sigma}^{(\text{post})} \right)^{-1} = \begin{pmatrix} x^2 + 1 & x \\ x & 2 \end{pmatrix} \Rightarrow \underline{\Sigma}^{(\text{post})} = \frac{1}{x^2 + 2} \begin{pmatrix} 2 & -x \\ -x & x^2 + 1 \end{pmatrix} \Rightarrow 2 \underline{\theta}^T \underline{A} \underline{\mu}$$

$$\left[ \underline{\Sigma}^{(\text{post})} \right] \underline{\mu}^{(\text{post})} = \begin{bmatrix} xy \\ y \end{bmatrix} \Rightarrow \underline{\mu}^{(\text{post})} = \frac{1}{x^2 + 2} \begin{bmatrix} yx \\ y \end{bmatrix} \leftarrow$$

$$p(m, c | y, x) = N \left( \begin{bmatrix} m \\ c \end{bmatrix} ; \underline{\mu}^{(\text{post})} \right) = \frac{1}{x^2 + 2} \begin{bmatrix} \frac{y+x}{2} \\ \frac{y-x}{2} \end{bmatrix}, \underline{\Sigma}^{(\text{post})} = \frac{1}{x^2 + 2} \begin{bmatrix} 2 & -x \\ -x & x^2 + 1 \end{bmatrix}$$

i)  $x = 0 \quad y = 0$



$$\underline{\Sigma}^{(\text{post})} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \underline{\mu}^{(\text{post})} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$y = mx + c + \varepsilon \sim N(0, 1)$$

$$y = m \cdot 0 + c + \varepsilon$$

$$0 = m \cdot 0 + c + \varepsilon$$

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

$$p(y|m, c) = N(y; \frac{m \cdot 0 + c}{\sqrt{1}}, 1)$$

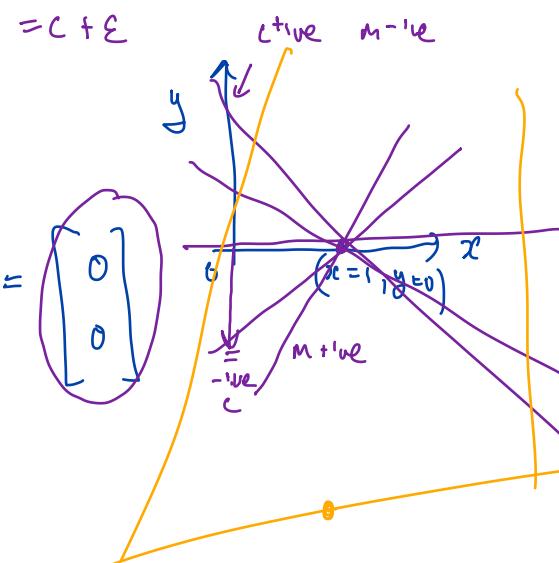
ii)  $x = 1 \quad y = 0$

$$\underline{\Sigma}^{(\text{post})} = \begin{bmatrix} 2/3 \\ -1/3 \\ -1/3 \\ 2/3 \end{bmatrix} \quad \underline{\mu}^{(\text{post})} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$y = mx + c + \varepsilon$$

$$0 = m + c + \varepsilon$$

$$0 = c + \varepsilon$$

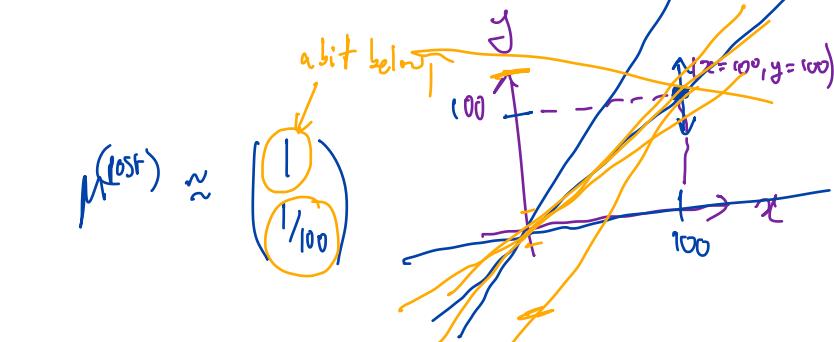


iii)  $x = 100 \quad y = 100$

$$\underline{\Sigma}^{(\text{post})} \approx \begin{bmatrix} 2/100^2 \\ -1/100 \\ 1 \end{bmatrix} \quad \underline{\mu}^{(\text{post})} \approx \begin{bmatrix} 1 \\ 1/100 \end{bmatrix}$$

$$y = mx + c + \varepsilon$$

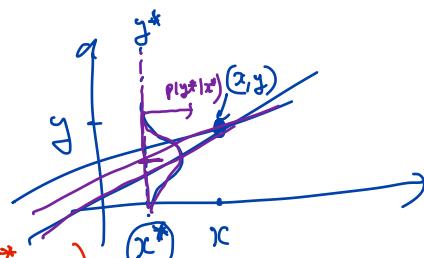
$$100 = m \cdot 100 + c + \varepsilon$$



How do we compute the predictive?

$$p(y^* | x^*, x, y) \leftarrow$$

$$y^* = mx^* + c + \epsilon^*$$



$$\textcircled{1} \quad p(y^* | x^*, x, y) = \int \underbrace{p(y^* | x^*, m, c)}_{\substack{\text{Gaussian} \\ \text{linear}}} \underbrace{p(m, c | x, y) dm dc}_{\substack{\text{Gaussian} \\ \text{prior} \\ \text{rule}}}$$

$$\int p(y^*, m, c | x^*, x, y) dm dc$$

$$p(y^* | x^*, x, y) = N(y^*; \mu_p^*, \sigma_p^2)$$

$$\mu^* = \mathbb{E}_{p(m, c | x, y) p(\epsilon^*)} \quad \text{[Red bracket under the first term]}$$

$$[y^*] = \mathbb{E}[m x^* + c + \epsilon^*]$$

$$\begin{cases} \textcircled{1} & m, c \sim p(m, c | x, y) = N([\bar{m}], \bar{\sigma}_m^2) \\ \textcircled{2} & y^* = mx^* + c + \epsilon^* \quad \epsilon^* \sim N(0, 1) \end{cases}$$

$$\mu^* = \mathbb{E}(m)x^* + \mathbb{E}(c) + \mathbb{E}(\epsilon^*)$$

$$\sigma_p^2 = \mathbb{E}_{p(m, c | x, y) p(\epsilon^*)} [ (y^* - \mu^*)^2 ]$$

$$\mu^* = \frac{y_x}{x^2+2} x^* + \frac{y}{x^2+2}$$

## Classification

### 10. 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  $\epsilon_n$  is Gaussian with mean 0 and variance  $\sigma^2$  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  $\sigma^2 \rightarrow \infty$ ?  $e^x =$

$$V^{(n)} = - \left| \frac{b_n - \underline{\mathbf{w}}^\top \underline{\mathbf{x}}_n}{\sigma^2} \right|$$

$$e^{-\frac{(V^{(n)})^2}{2}} = 1 + \frac{x^2}{2} + \frac{x^3}{3!} + \dots$$



## 11. 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) Interpreting the output of the softmax as  $y_i = p(t_i = 1 | \mathbf{W}, \mathbf{x})$  write down a cost-function for training this model 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?

