

70

University of Toronto  
STA 414/2104: Midterm , Winter 2015

Name: Jiayu Tang

Student ID: 1000930380

Read the following instructions carefully:

1. Do not turn the page until told to do so.
2. You may use a single side of  $8.5 \times 11$  inch aid sheet and a nonprogrammable calculator for this exam.
3. If a question asks you to do some calculations, you must *show your work* to receive full credit.
4. You can use either pen or pencil for the exam. But please be aware that you are not allowed to dispute any credit after the exam is returned if you use a pencil.
5. Lastly, enjoy the problems!!!

1. (10 mk) Show that the mode (i.e. the maximum) of the univariate Gaussian distribution

W

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

is given by  $\mu$ .

find  $x$ , to maximize  $N(x|\mu, \sigma^2)$ , is to maximize  $\exp(-\frac{1}{2\sigma^2}(x-\mu)^2)$ .

where  $\mu$  &  $\sigma$  are known

$\Rightarrow$  is equivalent to maximizing  $-\frac{1}{2\sigma^2}(x-\mu)^2$ . since  $h(a) = \exp(a)$  is an increasing function

$\Rightarrow$  that is to minimize  $(x-\mu)^2$ .

$(x-\mu)^2 \geq 0$ , it takes the equality holding only when  $x=\mu$ .

$\therefore$  When  $x=\mu$ ,  $N(x|\mu, \sigma^2)$  is maximized

$\therefore$  Mode of  $N(x|\mu, \sigma^2)$  is  $x=\mu$ .

2. (15 mk) Suppose that we have a dataset of observations:  $\mathbf{x} = (x_1, x_2, \dots, x_N)^\top$ , representing N observations of the scalar binary variable  $x$  (so that  $x_n \in \{0, 1\}$ ,  $n = 1, \dots, N$ ). Assume that the observations are i.i.d and are drawn from a Bernoulli distribution with parameter  $\mu$ .

- (a) (3 mk) Write down the log-likelihood for these  $N$  observations.

3 for each observation,  $\text{Bern}(x|1\mu) = \mu^x (1-\mu)^{1-x}$

Thus for  $N$  observations  $L(\mu; x_1, \dots, x_N) = \prod_{i=1}^N \mu^{x_i} (1-\mu)^{1-x_i}$  since  $x_i$ 's are i.i.d.

$$\log L(\mu; x_1, \dots, x_N) = \sum_{i=1}^N [x_i \log \mu + (1-x_i) \log (1-\mu)].$$

- (b) (5 mk) Maximizing the log-likelihood with respect to  $\mu$ , find the maximum likelihood

5 estimate  $\mu_{ML}$  of  $\mu$  (show your derivations).

$$\log L = \sum_{i=1}^N x_i \log \mu + \sum_{i=1}^N (1-x_i) \log (1-\mu)$$

$$= M \log \mu + (N-M) \log (1-\mu). \quad \text{where } M \text{ is the \# of observations} \\ \text{that } x_i = 1, \quad M = \sum_{i=1}^N x_i$$

$$\frac{\partial \log L}{\partial \mu} = \frac{M}{\mu} - \frac{N-M}{1-\mu} \stackrel{\text{set}}{=} 0.$$

$$\text{get } M_{ML} \cdot \frac{M}{N} = \frac{\sum_{i=1}^N x_i}{N}.$$

Checking the second derivative,  $\frac{\partial^2 \log L}{\partial \mu^2} = -\frac{M}{\mu^2} + \frac{N-M}{(1-\mu)^2} < 0$  at  $M_{ML}$

$$M_{ML} = \frac{M}{N}. \quad M \text{ is defined above.}$$

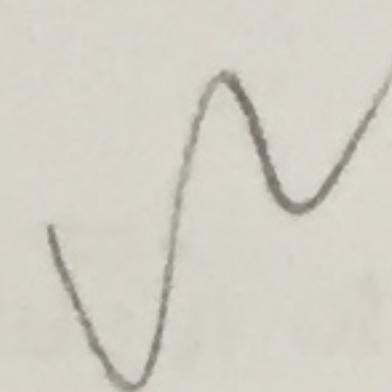
$$= \frac{\sum_{i=1}^N x_i}{N}.$$

(c) (7 mk) Show that the ML estimate of  $\mu$  is unbiased, i.e. show that:

7

$$\mathbb{E}[\mu_{ML}] = \mu$$

$$\begin{aligned}\mathbb{E}(\mu_{ML}) &= \mathbb{E}\left(\frac{\sum_{i=1}^N x_i}{N}\right) = \frac{1}{N} \mathbb{E}\left(\sum_{i=1}^N x_i\right) = \frac{1}{N} \left[ \sum_{i=1}^N \mathbb{E}(x_i) \right], \quad \mathbb{E}(x_i) = \mu \quad \forall i > 0 \\ &= \frac{1}{N} \cdot N \cdot \mu = \mu.\end{aligned}$$



3. (15 mk) Suppose that we have a dataset of observations:  $\mathbf{x} = (x_1, x_2, \dots, x_N)^\top$ , representing  $N$  observations of the scalar variable  $x$ . Assume that the observations are *i.i.d* and are drawn from a Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ .

- Suppose that the variance of a Gaussian is estimated using

$$\sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2$$

but with maximum likelihood estimate  $\mu_{ML}$  replaced with a true value  $\mu$  of the mean. Show that this estimator has the property that its expectation is given by the true variance  $\sigma^2$ , i.e. show that

$$\mathbb{E}[\sigma_{ML}^2] = \sigma^2.$$

$$\begin{aligned} \mathbb{E}(\sigma_{ML}^2) &= \mathbb{E}\left(\frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2\right) \text{ true.} \\ &= \frac{1}{N} \left[ \sum_{i=1}^N \mathbb{E}(x_i - \mu)^2 \right] \end{aligned}$$

$\mathbb{E}(x_i - \mu)^2 = \mathbb{E}[(x_i - \bar{x}_N)^2]$ , by the definition of variance,  $\mathbb{E}[(x_i - \mu)^2]$  is the true variance  $\sigma^2$  for  $x$ ,  $\mathbb{E}[x_i - \mu]^2 = \sigma^2 \quad \forall i$

Thus  $\mathbb{E}(\sigma_{ML}^2) = \frac{1}{N} \sum_{i=1}^N \sigma^2 = \sigma^2$ .

15

Remark

4. (20 mk) Consider a regression model  $y(\mathbf{x})$  with a squared loss function. In this case the expected loss can be written as:

$$\mathbb{E}[L] = \int \int (y(\mathbf{x}) - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

Show that the function  $y(\mathbf{x})$  for which this expected loss is minimized is given by the conditional expectation  $\mathbb{E}[t|\mathbf{x}] = \int t p(t|\mathbf{x}) dt$ .

$$\mathbb{E}[L] = \iint y(\mathbf{x})^2 \cdot p(\mathbf{x}, t) d\mathbf{x} dt + \iint t^2 p(\mathbf{x}, t) d\mathbf{x} dt - 2 \iint y(\mathbf{x}) \cdot t p(\mathbf{x}, t) d\mathbf{x} dt$$

$$= \int y(\mathbf{x})^2 \cdot p(\mathbf{x}) d\mathbf{x} + \int t^2 \cdot p(t) dt - 2 \int y(\mathbf{x}) d\mathbf{x} \int t p(t|\mathbf{x}) \cdot p(\mathbf{x}) dt$$

$$= \int y(\mathbf{x})^2 \cdot p(\mathbf{x}) d\mathbf{x} + \mathbb{E}[t^2] - 2 \int y(\mathbf{x}) p(\mathbf{x}) \cdot \mathbb{E}[t|\mathbf{x}] d\mathbf{x}$$

$$= \mathbb{E}[t^2] + \mathbb{E}_{\mathbf{x}}[y(\mathbf{x})^2 - 2y(\mathbf{x}) \mathbb{E}[t|\mathbf{x}]]$$

To minimize  $\mathbb{E}[L]$  is to minimize  $\mathbb{E}_{\mathbf{x}}[y(\mathbf{x})^2 - 2\mathbb{E}[t|\mathbf{x}] \cdot y(\mathbf{x})]$

$\forall \mathbf{x}$ ,  $y(\mathbf{x}) - 2\mathbb{E}[t|\mathbf{x}] \cdot y(\mathbf{x})$  is minimized at  $y(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}]$

Thus,  $\mathbb{E}_{\mathbf{x}}[y(\mathbf{x})^2 - 2\mathbb{E}[t|\mathbf{x}] \cdot y(\mathbf{x})] = \int [y(\mathbf{x})^2 - 2y(\mathbf{x}) \mathbb{E}[t|\mathbf{x}]] p(\mathbf{x}) d\mathbf{x}$  is always minimized at  $y(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}]$ .

$\therefore \mathbb{E}[L]$  is minimized at  $y(\mathbf{x}) = \int t p(t|\mathbf{x}) dt = \mathbb{E}[t|\mathbf{x}]$

Based on this,  $y(\mathbf{x}) = 0$  also works.  $\leftarrow$  correct ✓

5. (20 mk) Consider a linear basis function regression model. We observe a dataset of inputs  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  and corresponding real-valued target values  $\mathbf{t} = \{t_1, \dots, t_N\}$ . Assuming that the observations are *i.i.d*, the likelihood function can be written as:

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^\top \phi(\mathbf{x}_n), \sigma^2)$$

Consider a zero mean Gaussian prior governed by a single precision parameter  $\alpha$ :

evidence approximation.  $p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|0, \alpha^{-1}I)$

- (a) (3 mk) Write down the form of the marginal likelihood for this model.

1  $p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \alpha) = \int \prod_{n=1}^N \frac{1}{\sqrt{2\pi\alpha^2}} \exp\left(-\frac{1}{2\alpha^2} (\mathbf{x}_n^\top \mathbf{w})^2\right) p(\mathbf{w}) d\mathbf{w}$

$$p(\mathbf{t}|\mathbf{X}, \alpha, \beta) = \int p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w}|\alpha) d\mathbf{w} = (2\pi\alpha^2)^{-\frac{N}{2}} \exp\left\{-\frac{1}{2\alpha^2} \sum_{n=1}^N (\mathbf{x}_n^\top \mathbf{w})^2\right\}$$

$\sim N(\mathbf{t} | m_N, S_N)$

- 2 (b) (5 mk) Write down the form of predictive distribution of  $t$  for a new input  $\mathbf{x}^*$ .

Set  $\beta^{-1} = \alpha^2$ ,  $p(t^* | \mathbf{t}, \mathbf{x}^*, \mathbf{X}, \alpha, \beta) = N(t^* | m_N^\top \phi(\mathbf{x}^*), \alpha^2_N(\mathbf{x}^*))$

Bayesian approach where  $\alpha^2_N(\mathbf{x}^*) = \frac{1}{\beta} + \phi(\mathbf{x}^*)^\top S_N \phi(\mathbf{x}^*)$  → Indicate where this comes from.

$$S_N^{-1} = \alpha I + \beta \Phi^\top \Phi \quad m_N = \beta S_N \Phi^\top \mathbf{t}$$

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & & & \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

- (c) (12 mk) Show that maximizing the log of the posterior distribution with respect to  $w$  is equivalent to the minimization of the sum-of-squares error function with the addition of a quadratic regularization term, with regularization parameter  $\lambda = \alpha * \sigma^2$ .

12

the posterior fn  $p(w|D)$

since  $p(w|D) \propto p(w) \cdot p(D|w)$ .

$$p(w|D) = C \cdot \exp \left[ -\frac{\beta}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2 \right] \cdot \exp \left[ -\frac{\alpha}{2} w^\top w \right], \text{ where } C \text{ is a constant}$$

$$\log p(w|D) = -\frac{\beta}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2 - \frac{\alpha}{2} w^\top w + \text{constant}$$

To maximize  $\log p(w|D)$  w.r.t  $w$  is to minimize

$$-\frac{\beta}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2 - \frac{\alpha}{2} w^\top w$$

is to minimize

$$\beta \cdot \left[ \frac{1}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2 + \frac{1}{2} \frac{\alpha}{\beta} w^\top w \right], \quad \beta = \alpha^2 > 0.$$

Thus it is equivalent to minimize the error function.

$\frac{1}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2$  is the sum of the square error for the model

plus the regularization term.

$$\frac{1}{2} \lambda w^\top w \quad \text{where } \lambda = \frac{\alpha}{\beta} = \alpha \alpha^2.$$

6. (20 mk) Consider a Bayesian logistic regression model for a binary classification problem, parametrized by  $\mathbf{w}$ . Consider a dataset  $\{\mathbf{x}_n, t_n\}$ , where  $t_n \in \{0, 1\}$  and  $n = 1, \dots, N$ . Let us consider a zero mean Gaussian prior over  $\mathbf{w}$  governed by a single precision parameter  $\alpha$ :

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|0, \alpha^{-1}I) \quad (1)$$

- (a) (5 mk) Write down the posterior distribution over  $\mathbf{w}$  for the Bayesian logistic regression model. (Hint: posterior is proportional to the likelihood function times prior)

$$p(w|x,t) \propto p(t|x,w) \cdot p(w|\alpha). \quad p(w|\alpha) = \mathcal{N}(w|0, \alpha^{-1}I).$$

$$\begin{aligned} \text{Then } \log p(w|x,t) &= \log p(t|x,w) + \log p(w|\alpha) + \text{constant} \\ &= -\frac{1}{2} w^T \alpha I w + \sum_{n=1}^N [t_n \ln y_n + (1-t_n) \ln (1-y_n)] + \text{constant}. \end{aligned}$$

$$\text{Thus } p(w|x,t) \propto \exp\left(-\frac{1}{2}\alpha w^T w\right) \cdot \prod_{n=1}^N y_n^{t_n} (1-y_n)^{1-t_n} \checkmark$$

$$\text{where } y_n = \sigma(w^T x_n) = \frac{1}{1 + \exp(-w^T x_n)}.$$

The corresponding Gaussian approximation is. *Not needed*.

$$q(w) = \mathcal{N}(w|w_{\text{map}}, S_N^{-1}) : S_N^{-1} = \alpha I + \sum_n y_n (1-y_n) x_n x_n^T \quad w_{\text{map}} \text{ is the mode of the posterior}$$

- (b) (5 mk) Show that maximizing the log of the posterior distribution with respect to  $\mathbf{w}$  is equivalent to the minimization of the cross-entropy error function function with the addition of a quadratic regularization term (this corresponds to regularized cross-entropy error function).

$$\text{To maximize } \log p(w|x,t) = -\frac{\alpha}{2} w^T w + \sum_{n=1}^N \{t_n \ln y_n + (1-t_n) \ln (1-y_n)\} + \text{constant}$$

$$\text{is to minimize } -\sum_{n=1}^N \{t_n \ln y_n + (1-t_n) \ln (1-y_n)\} + \frac{\alpha}{2} w^T w.$$

$$= \sum_{n=1}^N E_n + \frac{\alpha}{2} w^T w$$

$$\text{where } E_n = -[t_n \ln y_n + (1-t_n) \ln (1-y_n)]$$

$$E_n = \bar{E}_n(y_n) = \begin{cases} -\ln y_n & \text{if } t_n = 1 \\ -\ln(1-y_n) & \text{if } t_n = 0. \end{cases} \checkmark$$

$\sum E_n$  is the cross-entropy error in total.

and  $\frac{\alpha}{2} w^T w$  is the corresponding regularization term with  $\lambda = \alpha$ .

- (c) (10 mk) Compute the derivative of the regularized cross-entropy function with respect to parameter vector  $\mathbf{w}$  (which is needed for parameter estimation).

W)

$$E[L] = -\sum_{n=1}^N [t_n \ln y_n + (1-t_n) \ln(1-y_n)] + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w}.$$

$$= \sum_{n=1}^N E_n + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \quad \checkmark$$

$$\frac{dE_n}{dy_n} = \frac{y_n - t_n}{y_n(1-y_n)} \quad \frac{dy_n}{dw} = y_n(1-y_n) \cdot x_n. \quad \text{since } \frac{dy_n}{dw} = \frac{d \sigma(w^T x_n)}{dw}$$

$$\therefore \frac{\partial E_n}{\partial w} = \frac{\partial E_n}{\partial y_n} \cdot \frac{\partial y_n}{\partial w} = (y_n - t_n) \cdot x_n.$$

$$\nabla_w [w^T w] = 2 \cdot w. \quad \checkmark$$

$$\text{Thus } \nabla_w E[L] = \sum_{n=1}^N (y_n - t_n) \cdot x_n + \alpha \cdot w.$$

95 / 100