

## 1 Probability Theory

**Problem 1 [2 points]** Your friend Mark Z. has a hepatic carcinoma checkup using a new medical technology. The result shows that Mark has the disease. The test is 98% accurate. In other words, this method detects hepatic carcinoma in 98 out of 100 times when it is there, and in 2/100 misses it; this method detects hepatic carcinoma in 2 of 100 cases when it is not there, and 98 of 100 times correctly returns a negative result. The previous research of hepatic carcinoma suggests that one in 2,000 people has this disease. What is the probability that Mark actually has hepatic carcinoma? Put down everything up to the point where you would use your calculator and stop there.

Assume  $T=1$  represent a positive test outcome,  $T=0$  represent a negative test outcome,  $D = 1$  mean Mark has the disease, and  $D = 0$  mean Mark doesn't have the disease. We have,  $P(T = 1|D = 1) = 0.98$ ,  $P(T = 0|D = 0) = 0.98$  and  $P(D = 1) = 0.005$ . The question is simplified to calculate  $P(D = 1|T = 1)$ . With Bayes' rule, we can derive:

$$\begin{aligned} P(D = 1|T = 1) &= \frac{P(T = 1|D = 1)P(D = 1)}{P(T = 1|D = 1)P(D = 1) + P(T = 1|D = 0)P(D = 0)} \\ &= \frac{0.98 \times 0.005}{0.98 \times 0.005 + 0.02 \times 0.995} = \end{aligned}$$

## 2 Linear Regression

**Problem 2 [2 points]** We have 1D input points  $\mathbf{X} = [0, 1, 2]$ , and corresponding 2D output  $\mathbf{Y} = \{[-1, 1], [1, -1], [2, -1]\}$ . We embed  $x_i$  into 2d with the basis function:

$$\Phi(0) = (1, 0)^T, \Phi(1) = (1, 1)^T, \Phi(2) = (2, 2)^T$$

The model becomes  $\hat{\mathbf{y}} = \mathbf{W}^T \Phi(x)$ . Compute the MLE for  $\mathbf{W}$ .

Put down everything up to the point where you would use your calculator and stop there.

$$\Phi^T = \begin{pmatrix} 1 & 1 & 2 \\ 0 & 1 & 2 \end{pmatrix}$$

$$\hat{\mathbf{W}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{Y} = \begin{pmatrix} -1 & 1 \\ 2 & -1.6 \end{pmatrix}$$

**Problem 3 [3 points]** Assume that  $\bar{\mathbf{x}} = 0$ . We have ridge regression, and we minimize:

$$J(\mathbf{w}) = (\mathbf{y} - \mathbf{X}\mathbf{w} - w_0\mathbf{1})^T (\mathbf{y} - \mathbf{X}\mathbf{w} - w_0\mathbf{1}) + \lambda \mathbf{w}^T \mathbf{w}$$

Derive the optimizer

$$\hat{\mathbf{w}}_{ridge} = (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\begin{aligned}
 J(\mathbf{w}) &= \mathbf{w}^T \mathbf{w} + \mathbf{w}^T (\mathbf{X}^T \mathbf{X}) \mathbf{w} - 2\mathbf{y}^T (\mathbf{X} \mathbf{w}) + \lambda \mathbf{x}^T \mathbf{x} - 2w_0 \mathbf{1}^T \mathbf{y} + 2w_0 \mathbf{1}^T \mathbf{X} \mathbf{w} + w_0 \mathbf{1}^T \mathbf{1} w_0 \\
 & \qquad \qquad \qquad w_0 \mathbf{1}^T \mathbf{X} \mathbf{w} = \bar{\mathbf{x}}^T \mathbf{w} = 0 \\
 \frac{\partial}{\partial \mathbf{w}} J(\mathbf{w}) &= [2\mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{X}^T \mathbf{y}] + 2\lambda \mathbf{w} = 0 \\
 \mathbf{w} &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}
 \end{aligned}$$

### 3 Logistic Regression

Consider the data in the following figure, where we fit the model  $p(y = 1 | \mathbf{x}, \mathbf{w}) = \sigma(w_0 + w_1 x_1 + w_2 x_2)$ . Suppose we fit the model by maximum likelihood, i.e., we minimize

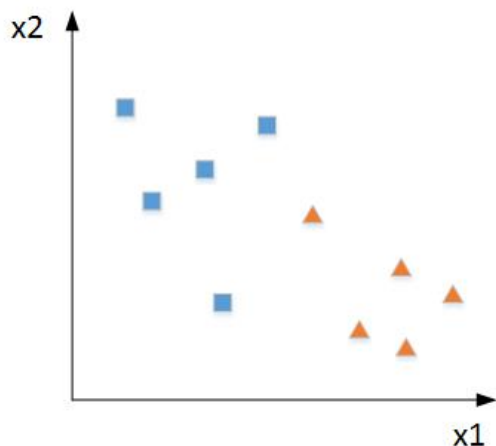
$$J(\mathbf{w}) = -L(\mathbf{w}, \mathcal{D}_{train})$$

where  $L(\mathbf{w}, \mathcal{D}_{train})$  is the log likelihood on the training set.

**Problem 4 [2 points]** Suppose we regularize only the  $w_0$  parameter, i.e., we minimize

$$J_0(\mathbf{w}) = -L(\mathbf{w}, \mathcal{D}_{train}) + \lambda w_0^2$$

and  $\lambda$  is a very large number. Sketch a possible decision boundary. Show your work.

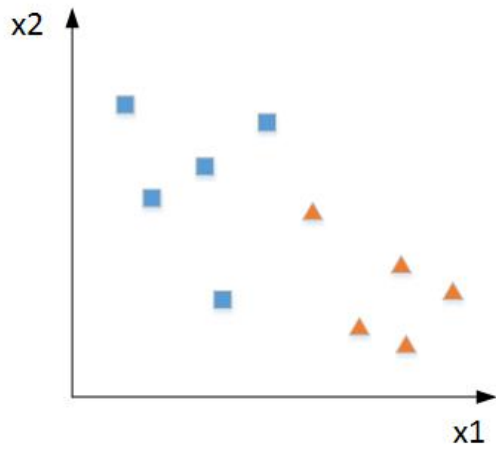


Heavily regularizing  $w_0$ , so the line must go through the origin. There are several possible line with different slopes.

**Problem 5 [2 points]** Now suppose we heavily regularize only the  $w_1$  parameter, i.e., we minimize

$$J_1(\mathbf{w}) = -L(\mathbf{w}, \mathcal{D}_{train}) + \lambda w_1^2$$

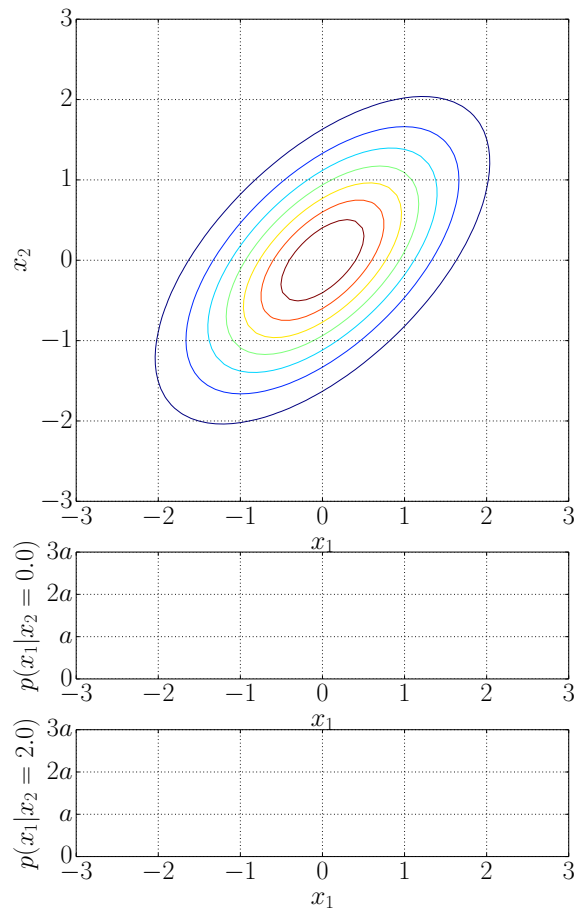
Sketch a possible decision boundary. Show your work.



Regularizing  $w_1$  makes the line horizontal since  $x_1$  is ignored.

## 4 Multivariate Gaussian

**Problem 6 [2 points]** The plot below shows a joint Gaussian distribution  $p(x_1, x_2)$ . Qualitatively (!) draw the conditionals  $p(x_1|x_2 = 0)$  and  $p(x_1|x_2 = 2)$  in the given coordinate systems. Note that the scaling  $a$  is arbitrary but fixed.



For a bivariate Gaussian distribution  $p(x_1, x_2) = \mathcal{N}(x_1, x_2 | \mu, \Sigma)$  with  $\mu = (\mu_1, \mu_2)^T$  and

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{pmatrix}$$

applying the formula for conditioning a Gaussian (e.g. Murphy, section 4.3.1) yields

$$p(x_1 | x_2) = \mathcal{N}(x_1 | \mu_{1|2}, \sigma_{1|2}) = \mathcal{N}\left(x_1 | \mu_1 + \frac{\sigma_{12}}{\sigma_2^2}(x_2 - \mu_2), \sigma_1^2 - \frac{\sigma_{12}^2}{\sigma_2^2}\right)$$

We see that while  $\mu_{1|2}$  depends on the value of  $x_2$ ,  $\sigma_{1|2}$  does not. That implies for the drawing:

- the shape of both conditional Gaussians is identical.
- $\mu_{1|2}$  can be roughly inferred graphically as the middle point between the intersections of the horizontal  $x_2$  lines with the iso-curves.

## 5 Kernels

We have  $\mathbf{x} = [x_1 \ x_2]^T$ . Given the mapping

$$\varphi(x) = [1 \ x_1^2 \ \sqrt{2}x_1x_2 \ x_2^2 \ \sqrt{2}x_1 \ \sqrt{2}x_2]^T$$

**Problem 7 [2 points]** Determine the kernel  $K(\mathbf{x}, \mathbf{y})$ . Simplify your answer.

$$K(\mathbf{x}, \mathbf{y}) = \varphi^T(\mathbf{x})\varphi(\mathbf{y}) = (1 + \mathbf{x}^T \mathbf{y})^2$$

## 6 Constrained optimization

Find the box with the maximum volume which has surface area no more than  $S \in \mathbb{R}^+$ .

**Problem 8 [2 points]** Derive the Lagrangian of the problem and the corresponding Lagrange dual function. Hint: set the parameters of the length, width and height to be  $l, w, h$  respectively.

**Problem 9 [3 points]** Solve the dual problem and give the solution to the original problem.

Surface area is  $S = 2lw + 2lh + 2hw$  and the volume is  $lwh$

The problem is equivalent to:

Minimize:  $f(l, w, h) = -lwh$

Subject to:  $h(l, w, h) = lw + lh + hw - \frac{S}{2} \leq 0$

$$L(l, w, h) = -lwh + \beta(lw + lh + hw - \frac{S}{2})$$

Computing the partial derivatives of  $L(l, w, h, \alpha)$  with respect to  $w, l, h$  gives

$$\begin{aligned}\frac{\partial L}{\partial l} &= -wh + \beta(w + h) = 0 \\ \frac{\partial L}{\partial w} &= -lh + \beta(l + h) = 0 \\ \frac{\partial L}{\partial h} &= -wl + \beta(w + l) = 0\end{aligned}$$

(1)

Solving this system of equations yields

$$l(\alpha) = w(\alpha) = h(\alpha) = 2\alpha.$$

Inserting this into  $L(l, w, h, \alpha)$  yields the Lagrange dual function

$$g(\alpha) = \min_{l, w, h} L(l, w, h, \alpha) = 4\alpha^3 - \alpha \frac{S}{2}.$$

Solving the dual problem:

$$0 = \frac{dg}{d\alpha}$$

subject to dual feasibility  $\alpha \leq 0$ , yields

$$\alpha = \left(\frac{S}{24}\right)^{1/2}$$

and thus

$$\begin{aligned}l = w = h &= \left(\frac{S}{6}\right)^{1/2} \\ \max(lwh) &= -\min(f) = \left(\frac{S}{6}\right)^{3/2}\end{aligned}$$

## 7 Neural Networks

This is an unfair mock question about material you have not yet seen.

**Problem 10 [2 points]** What is deep about learning? Describe in 2–3 sentences.

**Problem 11 [2 points]** What vanishes in the gradient?