

During the next few lectures we will be looking at the inference from training data problem as a *random* process modeled by the joint probability distribution over input (measurements) and output (say class labels) variables. In general, estimating the underlying distribution is a daunting and unwieldy task, but there are a number of constraints or "tricks of the trade" so to speak that under certain conditions make this task manageable and fairly effective.

To make things simple, we will assume a discrete world, i.e., that the values of our random variables take on a finite number of values. Consider for example two random variables X taking on k possible values x_1, \dots, x_k and H taking on two values h_1, h_2 . The values of X could stand for a Body Mass Index (BMI) measurement $weight/height^2$ of a person and H stands for the two possibilities h_1 standing for the "person being over-weight" and h_2 as the possibility "person of normal weight". Given a BMI measurement we would like to estimate the probability of the person being over-weight.

The joint probability $P(X, H)$ is a two dimensional array (2-way array) with $2k$ entries (cells). Each training example (x_i, h_j) falls into one of those cells, therefore $P(X = x_i, H = h_j) = P(x_i, h_j)$ holds the ratio between the number of hits into cell (i, j) and the total number of training examples (assuming the training data arrive i.i.d.). As a result $\sum_{ij} P(x_i, h_j) = 1$.

The projections of the array onto its vertical and horizontal axes by summing over columns or over rows is called *marginalization* and produces $P(h_j) = \sum_i P(x_i, h_j)$ the sum over the j 'th row is the probability $P(H = h_j)$, i.e., the probability of a person being over-weight (or not) before we see any measurement — these are called *priors*. Likewise, $P(x_i) = \sum_j P(x_i, h_j)$ is the probability $P(X = x_i)$ which is the probability of receiving such a BMI measurement to begin with — this is often called *evidence*. Note that, by definition, $\sum_j P(h_j) = \sum_i P(x_i) = 1$. In Fig. 2.1 we have that $P(h_1) = 14/22$, $P(h_2) = 8/22$ that is there is a higher prior probability of a person being over-weight than being of normal weight. Also $P(x_3) = 7/22$ is the highest meaning that we encounter BMI = x_3 with the highest probability.

The *conditional probability* $P(h_j | x_i) = P(x_i, h_j)/P(x_i)$ is the ratio between the number of hits in cell (i, j) and the number of hits in the i 'th column, i.e., the probability that the outcome is $H = h_j$ given the measurement $X = x_i$. In Fig. 2.1 we have $P(h_2 | x_3) = 3/7$. Note that

$$\sum_j P(h_j | x_i) = \sum_j \frac{P(x_i, h_j)}{P(x_i)} = \frac{1}{P(x_i)} \sum_j P(x_i, h_j) = P(x_i)/P(x_i) = 1.$$

Likewise, the conditional probability $P(x_i | h_j) = P(x_i, h_j)/P(h_j)$ is the number of hits in cell (i, j) normalized by the number of hits in the j 'th row and represents the probability of receiving BMI = x_i given the class label $H = h_j$ (over-weight or not) of the person. In Fig. 2.1 we have $P(x_3 | h_2) = 3/8$ which is the probability of receiving BMI = x_3 given that the person is known to be of normal weight. Note that $\sum_i P(x_i | h_j) = 1$.

h_1	2	5	4	2	1
h_2	0	0	3	3	2
	x_1	x_2	x_3	x_4	x_5

Figure 2.1: Joint probability $P(X, H)$ where X ranges over 5 discrete values and H over two values. Each entry contains the number of hits for the cell (x_i, h_j) . The joint probability $P(x_i, h_j)$ is the number of hits divided by the total number of hits (22). See text for more details.

The Bayes formula arises from:

$$P(x_i | h_j)P(h_j) = P(x_i, h_j) = P(h_j | x_i)P(x_i),$$

from which we get:

$$P(h_j | x_i) = \frac{P(x_i | h_j)P(h_j)}{P(x_i)}.$$

The left hand side $P(h_j | x_i)$ is called the *posterior* probability and $P(x_i | h_j)$ is called the *class conditional likelihood*. The Bayes formula provides a way to estimate the posterior probability from the prior, evidence and class likelihood. It is useful in cases where it is natural to compute (or collect data of) the class likelihood, yet it is not quite simple to compute directly the posterior. For example, given a measurement "12" we would like to estimate the probability that the measurement came from tossing a pair of dice or from spinning a roulette table. If $x = 12$ is our measurement, and h_1 stands for "pair of dice" and h_2 for "roulette" then it is natural to compute the class conditional: $P("12" | "pair of dice") = 1/36$ and $P("12" | "roulette") = 1/38$. Computing the posterior directly is much more difficult. As another example, consider medical diagnosis. Once it is known that a patient suffers from some disease h_j , it is natural to evaluate the probabilities $P(x_i | h_j)$ of the emerging symptoms x_i . As a result, in many inference problems it is natural to use the class conditionals as the basic building blocks and use the Bayes formula to invert those to obtain the posteriors.

The Bayes rule can often lead to unintuitive results — the one in particular is known as "base rate fallacy" which shows how a nonuniform prior can influence the mapping from likelihoods to posteriors. On an intuitive basis, people tend to ignore priors and equate likelihoods to posteriors. The following example is typical: consider the "Cancer test kit" problem² which has the following features: given that the subject has Cancer "C", the probability of the test kit producing a positive decision "+" is $P(+ | C) = 0.98$ (which means that $P(- | C) = 0.02$) and the probability of the kit producing a negative decision "-" given that the subject is healthy "H" is $P(- | H) = 0.97$ (which means also that $P(+ | H) = 0.03$). The prior probability of Cancer in the population is $P(C) = 0.01$. These numbers appear at first glance as quite reasonable, i.e., there is a probability of 98% that the test kit will produce the correct indication given that the subject has Cancer. What we are actually interested in is the probability that the subject has Cancer given that the test kit generated a positive decision, i.e., $P(C | +)$. Using Bayes rule:

$$P(C | +) = \frac{P(+ | C)P(C)}{P(+)} = \frac{P(+ | C)P(C)}{P(+ | C)P(C) + P(+ | H)P(H)} = 0.266$$

which means that there is a 26.6% chance that the subject has Cancer given that the test kit produced a positive response — by all means a very poor performance.

²This example is adopted from Yishai Mansour's class notes on Machine Learning.

If we draw the posteriors $P(h_1 | x)$ and $P(h_2 | x)$ using the probability distribution array in Fig. 2.1 we will see that $P(h_1 | x) > P(h_2 | x)$ for all values of X smaller than a value which is in between x_3 and x_4 . Therefore the decision which will minimize the probability of misclassification would be to choose the class with the maximal posterior:

$$h^* = \operatorname{argmax}_j P(h_j | x),$$

which is known as the Maximal A Posteriori (MAP) decision principle. Since $P(x)$ is simply a normalization factor, the MAP principle is equivalent to:

$$h^* = \operatorname{argmax}_j P(x | h_j)P(h_j).$$

In the case where information about the prior $P(h)$ is not known or it is known that the prior is uniform, then we obtain the Maximum Likelihood (ML) principle:

$$h^* = \operatorname{argmax}_j P(x | h_j).$$

The MAP principle is a particular case of a more general principle, known as "proper Bayes", where a *loss* is incorporated into the decision process. Let $l(h_i, h_j)$ be the loss incurred by deciding on class h_i when in fact h_j is the correct class. For example, the "0/1" loss function is:

$$l(h_i, h_j) = \begin{cases} 1 & i \neq j \\ 0 & i = j \end{cases}$$

The least-squares loss function is: $l(h_i, h_j) = \|h_i - h_j\|^2$ typically used when the outcomes are vectors in some high dimensional space rather than class labels. We define the *expected risk*:

$$R(h_i | x) = \sum_j l(h_i, h_j)P(h_j | x).$$

The proper Bayes decision policy is to minimize the expected risk:

$$h^* = \operatorname{argmin}_j R(h_j | x).$$

The MAP policy arises in the case $l(h_i, h_j)$ is the 0/1 loss function:

$$R(h_i | x) = \sum_{j \neq i} P(h_j | x) = 1 - P(h_i | x),$$

Thus,

$$\operatorname{argmin}_j R(h_j | x) = \operatorname{argmax}_j P(h_j | x).$$

2.1 Independence Constraints

At this point we may pause and ask what have we obtained? well, not much. Clearly, the inference problem is captured by the joint probability distribution and we do not need all these formulas to see this. How do we obtain the necessary data to fill in the probability distribution array to begin with? Clearly without additional simplifying constraints the task is not practical as the size of these kind of arrays are exponential in the number of variables. There are three families of simplifying constraints used in the literature:

- statistical independence constraints,
- parametric form of the class likelihood $P(x_i | h_j)$ where the inference becomes a density estimation problem,
- structural assumptions — latent (hidden) variables, graphical models.

Today we will focus on the first of these simplifying constraints — statistical independence properties.

Consider two random variables X and Y . The variables are statistically independent $X \perp Y$ if $P(X | Y) = P(X)$ meaning that information about the value of Y does not add anything about X . The independence condition is equivalent to the constraint: $P(X, Y) = P(X)P(Y)$. This can be easily proven: if $X \perp Y$ then $P(X, Y) = P(X | Y)P(Y) = P(X)P(Y)$. On the other hand, if $P(X, Y) = P(X)P(Y)$ then

$$P(X | Y) = \frac{P(X, Y)}{P(Y)} = \frac{P(X)P(Y)}{P(Y)} = P(X).$$

Let the values of X range over x_1, \dots, x_k and the values of Y range over y_1, \dots, y_l . The associated $k \times l$ 2-way array, $P(X = x_i, Y = y_j)$ is represented by the outer product $P(x_i, y_j) = P(x_i)P(y_j)$ of two vectors $P(X) = (P(x_1), \dots, P(x_k))$ and $P(Y) = (P(y_1), \dots, P(y_l))$. In other words, the 2-way array viewed as a *matrix* is of rank 1 and is determined by $k + l$ (minus 2 because the sum of each vector is 1) parameters rather than kl (minus 1) parameters.

Likewise, if $X_1 \perp X_2 \perp \dots \perp X_n$ are n statistically independent random variables where X_i ranges over k_i discrete and distinct values, then the n -way array $P(X_1, \dots, X_n) = P(X_1) \cdot \dots \cdot P(X_n)$ is an outer-product of n vectors and is therefore determined by $k_1 + \dots + k_n$ (minus n) parameters instead of $k_1 k_2 \dots k_n$ (minus 1) parameters³. Viewed as a tensor, the joint probability is a rank 1 tensor. The main point is that the statistical independence assumption reduced the representation of the multivariate joint distribution from *exponential to linear size*.

Since our variables are typically divided to measurement variables and an output/class variable H (or in general H_1, \dots, H_l), it is useful to introduce another, weaker form, of independence known as *conditional independence*. Variables X, Y are conditionally independent given H , denoted by $X \perp Y | H$, iff $P(X | Y, H) = P(X | H)$ meaning that given H , the value of Y does not add any information about X . This is equivalent to the condition $P(X, Y | H) = P(X | H)P(Y | H)$. The proof goes as follows:

- If $P(X | Y, H) = P(X | H)$, then

$$P(X, Y | H) = \frac{P(X, Y, H)}{P(H)} = \frac{P(X | Y, H)P(Y, H)}{P(H)} = \frac{P(X | Y, H)P(Y | H)P(H)}{P(H)} = P(X | H)P(Y | H)$$

- If $P(X, Y | H) = P(X | H)P(Y | H)$, then

$$P(X | Y, H) = \frac{P(X, Y, H)}{P(Y, H)} = \frac{P(X, Y | H)}{P(Y | H)} = P(X | H).$$

³I am a bit over simplifying things because we are ignoring here the fact that the entries of the array should be non-negative. This means that there are additional non-linear constraints which effectively reduce the number of parameters — but nevertheless it stays exponential.

Consider as an example, Joe and Mo live on opposite sides of the city. Joe goes to work by train and Mo by car. Let X be the event "Joe is late to work" and Y be the event "Mo is late for work". Clearly X and Y are not independent because there could be other factors. For example, a train strike will cause Joe to be late, but because of the strike there would be extra traffic (people using their car instead of the train) thus causing Mo to be late as well. Therefore, a third variable H standing for the event "train strike" would decouple X and Y .

From a computational standpoint, the conditional independence assumption has a similar effect to the unconditional independence. Let X range over k distinct values, Y range over r distinct values and H range over s distinct values. Then $P(X, Y, H)$ is a 3-way array of size $k \times r \times s$. Given that $X \perp Y \mid H$ means that $P(X, Y \mid H = h_i)$, a 2-way "slice" of the 3-way array along the H axis is represented by the outer-product of two vectors $P(X \mid H = h_i)P(Y \mid H = h_i)$. As a result the 3-way array is represented by $s(k + r - 2)$ parameters instead of $skr - 1$. Likewise, if $X_1 \perp \dots \perp X_n \mid H$ then the n -way array $P(X_1, \dots, X_n \mid H = h_i)$ (which is a slice along the H axis of the $(n + 1)$ -array $P(X_1, \dots, X_n, H)$) is represented by an outer-product of n vectors, i.e., by $k_1 + \dots + k_n - n$ parameters.

2.1.1 Example: Coin Toss

We will use the ML principle to estimate the bias of a coin. Let X be a random variable taking the value $\{0, 1\}$ and H would be our hypothesis taking a real value in $[0, 1]$ standing for the coin's bias. If the coin's bias is q then $P(X = 0 \mid H = q) = q$ and $P(X = 1 \mid H = q) = 1 - q$. We receive m i.i.d. examples x_1, \dots, x_m where $x_i \in \{0, 1\}$. We wish to determine the value of q . Given that $x_1 \perp \dots \perp x_m \mid H$, the ML problem we must solve is:

$$q^* = \operatorname{argmax}_q P(x_1, \dots, x_m \mid H = q) = \prod_{i=1}^m P(x_i \mid q) = \operatorname{argmax}_q \sum_i \log P(x_i \mid q).$$

Let $0 \leq \lambda \leq m$ stand for the number of '0' instances, i.e., $\lambda = |\{x_i = 0 \mid i = 1, \dots, m\}|$. Therefore our ML problem becomes:

$$q^* = \operatorname{argmax}_q \{\lambda \log q + (m - \lambda) \log(1 - q)\}$$

Taking the partial derivative with respect to q and setting it to zero:

$$\frac{\partial}{\partial q} [\lambda \log q + (m - \lambda) \log(1 - q)] = \frac{\lambda}{q} - \frac{m - \lambda}{1 - q} = 0,$$

produces the result:

$$q^* = \frac{\lambda}{m}.$$

2.1.2 Example: Gaussian Density Estimation

So far we considered constraints induced by conditional independent statements among the random variables as a means to reduce the space and time complexity of the multivariate distribution array. Another approach would be to assume some form of parametric form governing the entries of the array — the most popular assumption is Gaussian distribution $P(X_1, \dots, X_n) \sim N(\mu, E)$ with mean vector μ and covariance matrix E . The parameters of the density function are denoted by $\theta = (\mu, E)$ and for every vector $\mathbf{x} \in R^n$ we have:

$$P(\mathbf{x} \mid \theta) = \frac{1}{(2\pi)^{n/2} |E|^{1/2}} \exp^{-\frac{1}{2}(\mathbf{x} - \mu)^\top E^{-1}(\mathbf{x} - \mu)}.$$

Assume we are given an i.i.d sample of k points $S = \{\mathbf{x}_1, \dots, \mathbf{x}_k\}$, $\mathbf{x}_i \in R^n$, and we would like to find the Bayes optimal θ :

$$\theta^* = \underset{\theta}{\operatorname{argmax}} P(S | \theta),$$

by maximizing the likelihood (here we are assuming that the the priors $P(\theta)$ are equal, thus the maximum likelihood and the MAP would produce the same result). Because the sample was drawn i.i.d. we can assume that:

$$P(S | \theta) = \prod_{i=1}^k P(\mathbf{x}_i | \theta).$$

Let $L(\theta) = \log P(S | \theta) = \sum_i \log P(\mathbf{x}_i | \theta)$ and since Log is monotonously increasing we have that $\theta^* = \underset{\theta}{\operatorname{argmax}} L(\theta)$. The parameter estimation would be recovered by taking derivatives with respect to θ , i.e., $\nabla_{\theta} L = 0$. We have:

$$L(\theta) = -\frac{1}{2} \log |E| - \sum_{i=1}^k \frac{n}{2} \log(2\pi) - \sum_i \frac{1}{2} (\mathbf{x}_i - \mu)^\top E^{-1} (\mathbf{x}_i - \mu). \quad (2.1)$$

We will start with a simple scenario where $E = \sigma^2 I$, i.e., all the covariances are zero and all the variances are equal to σ^2 . Thus, $E^{-1} = \sigma^{-2} I$ and $|E| = \sigma^{2n}$. After substitution (and removal of items which do not depend on θ) we have:

$$L(\theta) = -nk \log \sigma - \frac{1}{2} \sum_i \frac{\|\mathbf{x}_i - \mu\|^2}{\sigma^2}.$$

The partial derivative with respect to μ :

$$\frac{\partial L}{\partial \mu} = \sigma^{-2} \sum_i (\mu - \mathbf{x}_i) = 0$$

from which we obtain:

$$\mu = \frac{1}{k} \sum_{i=1}^k \mathbf{x}_i.$$

The partial derivative with respect to σ is:

$$\frac{\partial L}{\partial \sigma} = \frac{nk}{\sigma} - \sigma^{-3} \sum_i \|\mathbf{x}_i - \mu\|^2 = 0,$$

from which we obtain:

$$\sigma^2 = \frac{1}{kn} \sum_{i=1}^k \|\mathbf{x}_i - \mu\|^2.$$

Note that the reason for dividing by n is due to the fact that $\sigma_1^2 = \dots = \sigma_n^2 = \sigma^2$, so that:

$$\frac{1}{k} \sum_{i=1}^k \|\mathbf{x}_i - \mu\|^2 = \sum_{j=1}^n \sigma_j^2 = n\sigma^2.$$

In the general case, E is a full rank symmetric matrix, then the derivative of eqn. (2.1) with respect to μ is:

$$\frac{\partial L}{\partial \mu} = E^{-1} \sum_i (\mu - \mathbf{x}_i) = 0,$$

and since E^{-1} is full rank we obtain $\mu = (1/k) \sum_i \mathbf{x}_i$. For the derivative with respect to E we note two auxiliary items:

$$\frac{\partial |E|}{\partial E} = |E|E^{-1}, \quad \frac{\partial}{\partial E} \text{trace}(AE^{-1}) = -(E^{-1}AE^{-1})^\top.$$

Using the fact that $\mathbf{x}^\top \mathbf{y} = \text{trace}(\mathbf{x}\mathbf{y}^\top)$ we can transform $\mathbf{z}^\top E^{-1} \mathbf{z}$ to $\text{trace}(\mathbf{z}\mathbf{z}^\top E^{-1})$ for any vector \mathbf{z} . Given that E^{-1} is symmetric, then:

$$\frac{\partial}{\partial E} \text{trace}(\mathbf{z}\mathbf{z}^\top E^{-1}) = -E^{-1} \mathbf{z}\mathbf{z}^\top E^{-1}.$$

Substituting $\mathbf{z} = \mathbf{x} - \mu$ we obtain:

$$\frac{\partial L}{\partial E} = -kE^{-1} + E^{-1} \left(\sum_i (\mathbf{x}_i - \mu)(\mathbf{x}_i - \mu)^\top \right) E^{-1} = 0,$$

from which we obtain:

$$E = \frac{1}{k} \sum_{i=1}^k (\mathbf{x}_i - \mu)(\mathbf{x}_i - \mu)^\top.$$

2.2 Incremental Bayes Classifier

Consider another application of conditional dependence which is the Bayes incremental rule. Suppose we have processed n examples $X^{(n)} = \{X_1, \dots, X_n\}$ and computed somehow $P(H | X^{(n)})$. We are given a new measurement X and wish to compute (update) the posterior $P(H | X^{(n)}, X)$. We will use the chain rule⁴:

$$P(X | Y, Z) = \frac{P(X, Y, Z)}{P(Y, Z)} = \frac{P(Z | X, Y)P(X | Y)P(Y)}{P(Z | Y)P(Y)} = \frac{P(Z | X, Y)P(X | Y)}{P(Z | Y)}$$

to obtain:

$$P(H | X^{(n)}, X) = \frac{P(X | X^{(n)}, H)P(H | X^{(n)})}{P(X | X^{(n)})}$$

from conditional independence, $P(X | X^{(n)}, H) = P(X | H)$. The term $P(X | X^{(n)})$ can be expanded as follows:

$$\begin{aligned} P(X | X^{(n)}) &= \sum_i \frac{P(X, X^{(n)} | H = h_i)P(H = h_i)}{P(X^{(n)})} = \sum_i \frac{P(X | H = h_i)P(X^{(n)} | H = h_i)P(H = h_i)}{P(X^{(n)})} \\ &= \sum_i P(X | H = h_i)P(H = h_i | X^{(n)}) \end{aligned}$$

After substitution we obtain:

$$P(H = h_i | X^{(n)}, X) = \frac{P(X | H = h_i)P(H = h_i | X^{(n)})}{\sum_j P(X | H = h_j)P(H = h_j | X^{(n)})}.$$

The old posterior $P(H | X^{(n)})$ is now the prior for the updated formula. Consider the following example⁵: We have a coin which could be either fair or biased towards Head at a probability of 0.6.

⁴this is based on the rule $P(X_1, \dots, X_n) = P(X_1 | X_2, \dots, X_n)P(X_2 | X_3, \dots, X_n) \cdots P(X_{n-1} | X_n)P(X_n)$

⁵adopted from Ron Rivest's 1994 class notes.

Let $H = h_1$ be the event that the coin is fair, and $H = h_2$ that the coin is biased. We start with prior probabilities $P(h_1) = 0.75$ and $P(h_2) = 0.25$ (we have a higher initial belief that the coin is fair). Suppose our first coin toss is a Head, i.e., $X_1 = "0"$. Then,

$$P(h_1 | x_1) = \frac{P(x_1 | h_1)P(h_1)}{P(x_1)} = \frac{0.5 * 0.75}{0.5 * 0.75 + 0.6 * 0.25} = 0.714$$

and $P(h_2 | x_1) = 0.286$. Our posterior belief that the coin is fair has gone down after a Head toss. Assume we have another measurement $X_2 = "0"$, then:

$$P(h_1 | x_1, x_2) = \frac{P(x_2 | h_1)P(h_1 | x_1)}{\text{normalization}} = \frac{0.5 * 0.714}{0.5 * 0.714 + 0.6 * 0.286} = 0.675,$$

and $P(h_2 | x_1, x_2) = 0.325$, thus our belief that the coin is fair continues to go down after Head tosses.

2.3 Bayes Classifier for 2-class Normal Distributions

For the last topic in this lecture consider the 2-class inference problem. We will encounter this problem in this course in the context of SVM and LDA. In the Bayes framework, if $H = \{h_1, h_2\}$ denotes the "class member" variable with two possible outcomes, then the MAP decision policy calls for making the decision based on data \mathbf{x} :

$$h^* = \operatorname{argmax}_{h_1, h_2} \{P(h_1 | \mathbf{x}), P(h_2 | \mathbf{x})\},$$

or in other words the class h_1 would be chosen if $P(h_1 | \mathbf{x}) > P(h_2 | \mathbf{x})$. The *decision surface* (as a function of \mathbf{x}) is therefore described by:

$$P(h_1 | \mathbf{x}) - P(h_2 | \mathbf{x}) = 0.$$

The questions we ask here is what would the Bayes optimal decision surface be like if we assume that the two classes are normally distributed with different means and the same covariance matrix? What we will see is that under the condition of equal priors $P(h_1) = P(h_2)$ the decision surface is a hyperplane — and not only that, it is the same hyperplane produced by LDA.

Claim 1 *If $P(h_1) = P(h_2)$ and $P(\mathbf{x} | h_1) \sim N(\mu_1, E)$ and $P(\mathbf{x} | h_2) \sim N(\mu_2, E)$, then the Bayes optimal decision surface is a hyperplane $\mathbf{w}^\top (\mathbf{x} - \mu) = 0$ where $\mu = (\mu_1 + \mu_2)/2$ and $\mathbf{w} = E^{-1}(\mu_1 - \mu_2)$. In other words, the decision surface is described by:*

$$\mathbf{x}^\top E^{-1}(\mu_1 - \mu_2) - \frac{1}{2}(\mu_1 + \mu_2)^\top E^{-1}(\mu_1 - \mu_2) = 0. \quad (2.2)$$

Proof: The decision surface is described by $P(h_1 | \mathbf{x}) - P(h_2 | \mathbf{x}) = 0$ which is equivalent to the statement that the ratio of the posteriors is 1, or equivalently that the log of the ratio is zero, and using Bayes formula we obtain:

$$0 = \log \frac{P(\mathbf{x} | h_1)P(h_1)}{P(\mathbf{x} | h_2)P(h_2)} = \log \frac{P(\mathbf{x} | h_1)}{P(\mathbf{x} | h_2)}.$$

In other words, the decision surface is described by

$$\log P(\mathbf{x} | h_1) - \log P(\mathbf{x} | h_2) = -\frac{1}{2}(\mathbf{x} - \mu_1)^\top E^{-1}(\mathbf{x} - \mu_1) + \frac{1}{2}(\mathbf{x} - \mu_2)^\top E^{-1}(\mathbf{x} - \mu_2) = 0.$$

After expanding the two terms we obtain eqn. (2.2). \square