

Bayesian Estimation (DHS 3.3)

Bayesian Classifier (DHS 3.3.1)

$P(S_k|\underline{x}) \geq P(S_j|\underline{x})$ for all $j=1$ to $k \Rightarrow \underline{x} \in S_k$

$$P(S_k|\underline{x}) \propto P(\underline{x}|S_k)P(S_k)$$

What do we do if $P(S_j)$ and $P(\underline{x}|S_j)$ are unknown? Compute them using all information we have.

Given \underline{z} , the set of samples, compute the posterior probabilities $P(S_k|\underline{x},\underline{z})$ ← Final Goal

(i.e., use the training samples to compute the class-conditional density and prior density)

From Bayes theorem, $P(S_k|\underline{x},\underline{z}) = p(\underline{x}|S_k,\underline{z})P(S_k|\underline{z}) / \sum_{j=1}^K p(\underline{x}|S_j,\underline{z}) P(S_j|\underline{z})$

Assume $P(S_j|\underline{z}) = P(S_j)$ and $P(S_j)$ are known

Subdivide \underline{z} : $\underline{z}_1, \underline{z}_2, \dots, \underline{z}_k$

where \underline{z}_i contains all prototypes in class S_i

Assume $p(\underline{x}|S_i,\underline{z}) = p(\underline{x}|S_i,\underline{z}_i)$

$$P(S_k|\underline{x},\underline{z}) = p(\underline{x}|S_k, \underline{z}_k) P(S_k) / \sum_{j=1}^K p(\underline{x}|S_j,\underline{z}_j) P(S_j)$$

That is treat each class separately,

$P(\underline{x}|S_k)$ has known parametric form $\Rightarrow p(\underline{x}|S_k, \underline{\theta})$ is known.

The goal is to determine the likelihood, $p(\underline{x}|S_k, \underline{z}_k)$ using the prototypes \underline{z}_k . ← Goal

Parameter Distribution (DHS 3.3.2)

Now, our goal is to compute $p(x|z_k)$, which can be computed from

$$p(x|z_k) = \int p(x, \underline{\theta} | z_k) d\underline{\theta}$$

since the selection of x and the selection of the training samples in z_k is done independently, rewrite this

$$p(x|z_k) = \int p(x|\underline{\theta})p(\underline{\theta} | z_k)d\underline{\theta} \quad (*)$$

This equation links the class-conditional density $p(x|z_k)$ to the posterior density $p(\underline{\theta} | z_k)$ for the unknown parameter vector.

If $p(\underline{\theta} | z_k)$ peaks very sharply about the some value $\hat{\underline{\theta}}$, we obtain $p(x|z_k) \approx p(x|\hat{\underline{\theta}})$.

This means the result is obtained by substituting $\hat{\underline{\theta}}$ (the estimate) for the true parameter.

Bayesian Parameter Estimation: Gaussian Case (DHS 3.4)

Goal: compute **(Goal 1)** $p(\underline{\theta} | z)$ and then **(Goal 2)** $p(x | z)$ in (*). ← Two Goals

(Goal 1) Learning the mean of a normal density: get $p(\mu|z_k)$ (DHS 3.4.1): Univariate Case

Assumption: $p(x|\mu) = N(x, \mu, \sigma^2)$

μ =unknown

σ^2 =known

$\theta = \mu$

Objective: determine $p(\mu|z_k)$

Assume: $p(\mu) = N(\mu, \mu_0, \sigma_0^2)$ known

μ_0 =known

= a priori guess of parameter μ

σ_0^2 =known

= variance or uncertainty of the guess μ_0

Bayes theorem $p(\mu|z) = [p(z|\mu)p(\mu)] / [\int p(z|\mu) p(\mu) d\mu]$

Denominator = $p(z) = 1/\alpha$

$$\{z_k\} = \{x_1, x_2, \dots, x_J\}$$

= set of prototypes from same class k , independently drawn from the population

Drop $k \rightarrow z$ & Drop S_k

x_i 's independently drawn $\Rightarrow p(z|\mu) = \prod_{j=1}^J p(x_j|\mu)$

$$p(\mu|z) = \alpha \prod_{j=1}^J p(x_j|\mu)p(\mu)$$

$$= \alpha p(\mu) \prod_{j=1}^J p(x_j|\mu)$$

$p(\mu) = N(\mu, \mu_0, \sigma_0^2)$ (known)

$p(x_j|\mu) = N(x_j, \mu, \sigma^2)$ (μ unknown and σ^2 known)

$$p(\mu|z) = \alpha \frac{1}{\sqrt{2\pi}\sigma_0} \exp\left\{-\frac{1}{2}\left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right\} \prod_{j=1}^J \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2}\left(\frac{x_j - \mu}{\sigma}\right)^2\right\}$$

$$= \alpha' \exp\left\{-\frac{1}{2}\left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right\} \exp\left\{-\frac{1}{2}\sum_{j=1}^J \left(\frac{x_j - \mu}{\sigma}\right)^2\right\}$$

$$= \alpha' \exp\left\{-\frac{1}{2\sigma_0^2}\mu^2 + \frac{\mu_0}{\sigma_0^2}\mu - \frac{1}{2}\frac{\mu_0^2}{\sigma_0^2}\right\} \exp\left\{-\frac{1}{2\sigma^2}\sum_{j=1}^J x_j^2 + \frac{\mu}{\sigma^2}\sum_{j=1}^J x_j - \frac{1}{2\sigma^2}J\mu^2\right\}$$

$$= \alpha'' \exp\left\{-\frac{1}{2}\left[\left(\frac{1}{\sigma_0^2} + \frac{J}{\sigma^2}\right)\mu^2 - 2\left(\frac{\mu_0}{\sigma_0^2} + \frac{Jm_J}{\sigma^2}\right)\mu\right]\right\} \quad (**)$$

$$\text{where } m_J = \frac{1}{J}\sum_{j=1}^J x_j$$

Since $p(\mu|z)$ =normal

$$\text{Now set, } p(\mu|z) = \frac{1}{\sqrt{2\pi}\sigma_J} \exp\left\{-\frac{1}{2}\left(\frac{\mu - \mu_J}{\sigma_J}\right)^2\right\} = N(\mu, \mu_J, \sigma_J^2)$$

$$= \beta \exp\left\{-\frac{1}{2}\left[\frac{1}{\sigma_J^2}\mu^2 - \frac{2\mu_J}{\sigma_J^2}\mu + \frac{\mu_J^2}{\sigma_J^2}\right]\right\} \quad \text{and compare to (**)}$$

$$\Rightarrow \frac{1}{\sigma_J^2} = \frac{1}{\sigma_0^2} + \frac{J}{\sigma^2}$$

$$\frac{\mu_J}{\sigma_J^2} = \frac{\mu_0}{\sigma_0^2} + \frac{Jm_J}{\sigma^2}$$

$$\text{So } \sigma_J^2 = \frac{\sigma^2 \sigma_0^2}{\sigma^2 + J\sigma_0^2}$$

$$\mu_J = \sigma_J^2 \left[\frac{\mu_0}{\sigma_0^2} + \frac{Jm_J}{\sigma^2} \right]$$

Therefore $p(\mu|z)=N(\mu, \mu_J, \sigma_J^2)$

$$\mu_J = \frac{\sigma^2}{\sigma^2 + J\sigma_0^2} \mu_0 + \frac{J\sigma_0^2}{\sigma^2 + J\sigma_0^2} m_J$$

$$\sigma_J^2 = \frac{\sigma^2 \sigma_0^2}{\sigma^2 + J\sigma_0^2}, \quad m_J = \frac{1}{J} \sum_{j=1}^J x_j$$

These equations show how prior is combined with empirical information in samples to obtain a posteriori density $p(\underline{\theta} | \underline{z})$

μ_J is best estimate of μ after J observations.

σ_J is the uncertainty in the estimate μ_J .

This behavior is called Bayesian Learning (note DHS Fig. 3.2)

(Goal 2) Now get $p(\underline{x}|z)$: again univariate case (Read DHS 3.4.2)

Now we know $p(\mu|z)$.

Need to obtain the class-conditional density, $p(x|z)=p(x|S_k, z_k)$

Recall (*) $\Rightarrow p(x|z)=\int p(x|\mu)p(\mu|z)d\mu$

$$p(x|z)=\int \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right\} \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left\{-\frac{1}{2}\left(\frac{\mu-\mu_j}{\sigma_j}\right)^2\right\} d\mu$$

$$p(x|z)= \frac{1}{2\pi\sigma\sigma_j} \exp\left\{-\frac{1}{2}\left(\frac{(x-\mu_j)^2}{\sigma^2 + \sigma_j^2}\right)\right\} f(\sigma, \sigma_j)$$

$$\text{where } f(\sigma, \sigma_j) = \int \exp\left\{-\frac{1}{2} \frac{\sigma^2 + \sigma_j^2}{\sigma^2 \sigma_j^2} \left(\mu - \frac{\sigma_j^2 \mu + \sigma^2 \mu_j}{\sigma^2 + \sigma_j^2}\right)^2\right\} d\mu$$

This makes $p(x|z)$ as a function of x is proportional to $\exp\left\{-\frac{1}{2}\left(\frac{(x-\mu_j)^2}{\sigma^2 + \sigma_j^2}\right)\right\}$

Therefore $p(x|z)$ is normally distributed with mean μ_j and $\sigma^2 + \sigma_j^2$:

$$\therefore p(x|z)=N(x, \mu_j, \sigma^2 + \sigma_j^2)$$

[Summary] Bayesian Estimation Procedure

1. Estimate μ_j and σ_j by the boxed formulas and substitute into the equation for $p(\mu|z)$.
2. Find $p(x|z)$ from above.
3. Use the following to find, $P(S_k|\underline{x}, z)$, given $p(x|z)=p(x|S, z)$ and $p(S)$

$$p(S_k|\underline{x}, z) = p(\underline{x}|S_k, z) p(S_k) / \sum_{j=1}^K p(\underline{x}|S_j, z) p(S_j)$$

(for class S_k , $p(x|z)=p(x|S_k, z_k)$)

Multivariate Extension

$$p(x | \mu) = N(x, \mu, \Sigma)$$

$$p(\mu) = N(\mu, \mu_0, \Sigma_0)$$

$$p(x | z) = N(x, \mu_J, \Sigma + \Sigma_J)$$

μ_0, Σ_0, Σ are known. μ is unknown

where

$$\begin{aligned} \mu_J &= \Sigma_0 [\Sigma_0 + \Sigma/J]^{-1} m_J + (1/J) \Sigma [\Sigma_0 + (1/J) \Sigma]^{-1} \mu_0 \\ \Sigma_J &= \Sigma_0 [\Sigma_0 + 1/J \Sigma]^{-1} (\Sigma/J) \end{aligned}$$

$$m_J = \frac{1}{J} \sum_{j=1}^J x_j$$

μ_0 = initial guess of μ

Σ_0 = initial uncertainty.

Generalize the above technique => General Bayesian Learning

(Again, do for each class S_k separately)

1. General form of $p(\underline{x}|\theta)$ is known, but θ is not known exactly.
2. Initial knowledge about $\underline{\theta}$ is available as $p(\underline{\theta})$. Rest of our knowledge is a set of samples $\{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_J\} = \underline{z}$ of known classification drawn from a population of unknown density $p(\underline{x})$. (\underline{x}_j is independent of \underline{x}_k)

Objective: Compute $p(\underline{x}|\underline{z})$ to obtain $P(S_k|\underline{x}, \underline{z})$

Procedure

1. $p(\underline{x}|\underline{z}) = \int p(\underline{x}|\theta)p(\theta|\underline{z})d\theta$ (*)

where

2. $p(\theta|\underline{z}) = p(\underline{z}|\theta)p(\theta) / \int p(\underline{z}|\theta)p(\theta)d\theta$

where

3. $p(\underline{z}|\theta) = \prod_{j=1}^J p(\underline{x}_j|\theta)$

Note: If $p(\underline{z}|\theta)$ peaks at $\underline{\theta} = \hat{\underline{\theta}}$, then $p(\theta|\underline{z})$ will peak at $\underline{\theta} = \hat{\underline{\theta}}$

ML: θ that maximizes $p(\underline{z}|\theta)$.

Recursive Bayesian Learning

From 3, $p(z^{(j)}|\theta) = p(x_j|\theta) \prod_{j=1}^{J-1} p(x_j|\theta)$

or $p(z^{(j)}|\theta) = p(x_j|\theta) p(z^{(j-1)}|\theta)$

From 2, $p(\theta|z^{(j)}) = \frac{[p(x_j | \theta)p(z^{(j-1)} | \theta)]p(\theta)}{\int [p(x_j | \theta)p(z^{(j-1)} | \theta)]p(\theta)d\theta}$

Use Bayes rule

$p(z^{(j-1)} | \theta)p(\theta) = p(\theta | z^{(j-1)})p(z^{(j-1)})$

$p(\theta | z^{(j)}) = \frac{p(x_j | \theta)p(\theta | z^{(j-1)})p(z^{(j-1)})}{\int p(x_j | \theta)p(\theta | z^{(j-1)})p(z^{(j-1)})d\theta}$

$p(\theta|z^{(0)}) = p(\theta)$

Sequence: $p(\theta)$ (initial guess with no data)

$p(\theta|z^{(1)}) = p(\theta|x_1)$

$p(\theta|z^{(2)}) = p(\theta|x_1,x_2)$

...

Usually converges to a delta function.

Which Method is Better? Maximum-Likelihood or Bayes Methods (DHS 3.5.1)

- Computational Complexity: ML is preferred since it requires merely differential calculus techniques or gradient search w.r.t. parameters. Bayes methods require complex multidimensional integration.
- ML will be easier to interpret and understand, but Bayesian gives a weighted average of models or parameters, leading to solutions more complicated and harder to understand.
- Bayesian uses more information brought into the problem than ML.
- If more reliable information available, Bayes gives better results.
- Bayesian with a flat prior gives same results as ML.
- Bayesian balances between the accuracy of the estimation and its variance.
- Three Error Sources
 - Bayes error: due to overlapping densities. Cannot be eliminated.
 - Model error: due to an incorrect model. With better models, can be reduced.
 - Estimation error: due to a finite sample. Reduced with more training data.

Problems of Dimensionality (DHS 3.7)

- How classification accuracy depends upon the dimensionality and the amount of training data
- The computational complexity of designing the classifier.
- For Bayes classifier, the most useful features are the ones that offer bigger differences between the means than the standard deviations, thus reducing the probability of error.
- An obvious way is to introduce new independent features.
- If performance of a classifier is poor, it is natural to utilize new features, particularly ones that will help separate the class pairs most frequently confused.
- But increasing the number of features increases the cost and complexity of both the feature extractor and the classifier.
- In general, the performance should improve
- See DHS Fig. 3.3

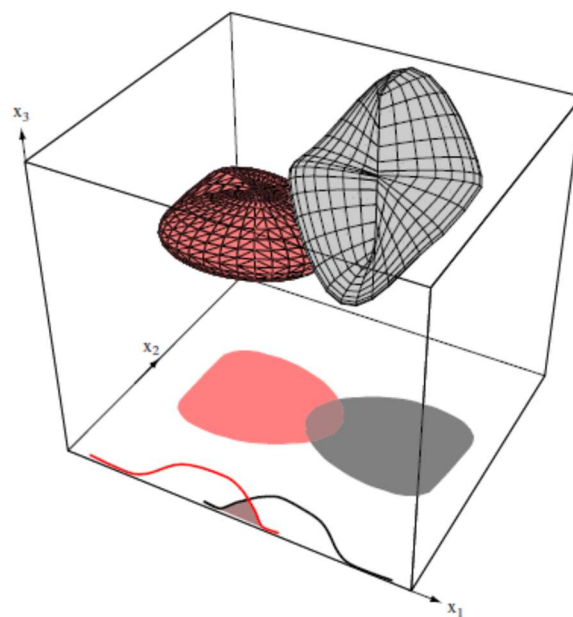


Figure 3.3: Two three-dimensional distributions have nonoverlapping densities, and thus in three dimensions the Bayes error vanishes. When projected to a subspace — here, the two-dimensional $x_1 - x_2$ subspace or a one-dimensional x_1 subspace — there can be greater overlap of the projected distributions, and hence greater Bayes errors.

Overfitting (DHS 3.7.3)

- See DHS Fig. 3.4

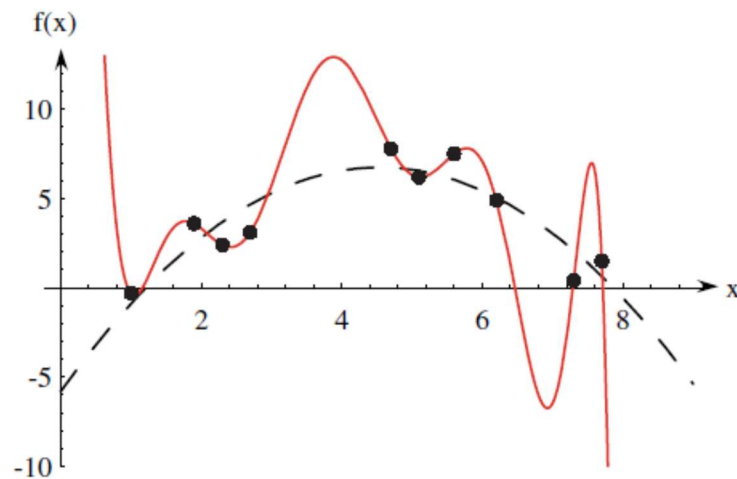


Figure 3.4: The “training data” (black dots) were selected from a quadratic function plus Gaussian noise, i.e., $f(x) = ax^2 + bx + c + \epsilon$ where $p(\epsilon) \sim N(0, \sigma^2)$. The 10th degree polynomial shown fits the data perfectly, but we desire instead the second-order function $f(x)$, since it would lead to better predictions for new samples.

** Principle Component Analysis (DHS 3.8.1) and Fisher Linear Discriminant (DHS 3.8.2) will be covered later in Unsupervised Classification.*