Estimation of Unknown Prob.
Density Functions

− Lecture 4 −

- Until now, we assumed that pdf's are known
- This is not the common case:
  ↳ in many problems, the underlying pdf has to be estimated
    from the available data.

- suppose that we can reasonably assume $p(x|w_i)$ is a Normal Density
with mean $\mu_i$ and cov. matrix $\Sigma_i$ (although we don't know the exact
values of these quantities)
  ↳ the problem is simplified then from estimating an unknown
function $p(x|w_i)$ to the one of estimating the parameters $\mu_i$ and $\Sigma_i$

# Maximum Likelihood Parameter Estimation

- views the parameters as quantities whose values are fixed but unknown!

- maximizes the ==probability of obtaining== the samples actually observed.

Suppose we have a collection of samples from c classes:

$$D = \{D_1, D_2, \ldots D_c\} \qquad : \text{dataset in } c \text{ classes}$$

- samples from $D_j$ have been drawn independently

according to $p(\vec{x} | w_j)$

such samples are i.i.d → independent and identically distri.

Assume $p(\vec{x} \mid w_j)$ has a known parametric form,
determined uniquely by $\vec{\theta_j}$

ex: $p(\vec{x} \mid w_j) \sim N(\mu_j, \xi_j) \quad \rightarrow \vec{\theta_j} = \{\mu_j, \xi_j\}$

to show the dependence of $p(\vec{x} \mid w_j)$ on $\theta_j$:

$$p(\vec{x} \mid w_j, \theta_j)$$

Aim: use the information provided by the training samples to
obtain good estimates for the unknown vectors:

$$\{\theta_1, \theta_2, \ldots \theta_c\}$$

Now, assume that samples $D_i$ give no information about $\theta_j$.

$\hookrightarrow$ Hence, parameters for the different classes are functionally independent.

$\hookrightarrow$ So, we can work with each class separately:

call $\theta$ to params of a class (not using subscript $\theta_j$ anymore)

Suppose $D$ contains $n$ samples: $\{x_1, x_2, \dots x_n\}$. Since samples are i.i.d:

$$p(D|\theta) = \prod_{k=1}^{n} p(x_k|\theta)$$

likelihood function of $\theta$ w.r.t the set of samples $D$

Find the value $\hat{\theta}$, that maximizes $p(D|\theta)$

Find $\hat{\theta}$ that maximizes $p(D|\theta)$ :

For analytic simplicity, use logarithm of the likelihood

$\hookrightarrow$ monotonically increasing function : no problem.

if the # of params is $p$: $\vec{\theta} = [\theta_1, \theta_2, \ldots \theta_p]^t$

$\nabla \to$ gradient op. $\qquad \nabla_\theta = \left[\frac{\partial}{\partial\theta_1}, \frac{\partial}{\partial\theta_2}, \cdot\frac{\partial}{\partial\theta_p}\right]^t$

$$l(\theta) = \ln p(D|\theta) \qquad = \ln \prod_{k=1}^{n} p(x_k|\theta)$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \, l(\theta)$$

$$l(\theta) = \sum_{k=1}^{n} \ln p(x_k|\theta)$$

$$\ell(\theta) = \sum_{k=1}^{\hat{n}} \ln p(x_k | \theta)$$

$$\nabla_\theta \ell(\theta) = \sum_{k=1}^{\hat{n}} \nabla_\theta \ln p(x_k | \theta)$$

$$\nabla_\theta \ell(\theta) = 0 \quad \rightarrow \text{set of } p \text{ equations.}$$

- solution to $\hat{\theta}$ could represent a true global max, or $\left\{ \begin{array}{l} \text{check each} \\ \text{sln. individually.} \end{array} \right.$

  a local max

## Example: The Gaussian Case: Unknown $\mu$!

Assume that samples are drawn from a multivariate normal population with mean $\mu$ and cov. $\Sigma$. Assume only $\mu$ unknown!

$$\ln p(x_k | \mu) = -\frac{1}{2} \ln [(2\pi)^d |\Sigma|] - \frac{1}{2}(x_k - \mu)^t \Sigma^{-1}(x_k - \mu)$$

$$\nabla_\mu \ln p(x_k | \mu) = \Sigma^{-1}(x_k - \mu)$$

the maximum likelihood est. of $\mu$ must satisfy:

$$\sum_{k=1}^{\hat{n}} \Sigma^{-1}(x_k - \hat{\mu}) = 0 \qquad\qquad n.\hat{\mu} = \sum_{k=1}^{\hat{n}} x_k$$

$$\sum_{k=1}^{\hat{n}} (x_k - \hat{\mu}) = 0$$

$$\sum_{k=1}^{\hat{n}} \hat{\mu} = \sum_{k=1}^{\hat{n}} x_k$$

$$\boxed{\hat{\mu} = \frac{1}{n} \sum_{k=1}^{\hat{n}} x_k}$$

arithmetic avg. of the observed samples

# Example 2: The Gaussian Case: Unknown $\mu$ and $\sigma^2$ → Univariate Normal Density

$$\theta = \{\mu, \sigma^2\} \quad \rightarrow \text{parameters to be estimated}$$

$$\underset{\theta_1 \quad \theta_2}{\downarrow \quad \downarrow}$$

$$\ln p(x_k | \theta) = -\frac{1}{2} \ln 2\pi \theta_2 - \frac{1}{2\theta_2}(x_k - \theta_1)^2$$

$$\nabla_\theta l = \nabla_\theta \ln p(x_k|\theta) = \begin{bmatrix} \frac{1}{\theta_2}(x_k - \theta_1) \\ \\ -\frac{1}{2\theta_2} + \frac{(x_k - \theta_1)^2}{2\theta_2^2} \end{bmatrix}$$

$$\nabla_\theta l = 0 \qquad (\text{max. likelihood cond})$$

$$\nabla_\theta \, l = \nabla_\theta \, \ln p(x_k | \theta) = \begin{bmatrix} \dfrac{1}{\theta_2} \quad (x_k - \theta_1) \\[4mm] -\dfrac{1}{2\theta_2} + \dfrac{(x_k - \theta_1)^2}{2\theta_2^2} \end{bmatrix}$$

$$\sum_{k=1}^{n} \frac{1}{\hat{\theta}_2} (x_k - \hat{\theta}_1) = 0 \qquad (1) \implies \hat{\theta}_1 = \hat{\mu} = \frac{1}{n} \sum_{k=1}^{n} x_k$$

$$-\sum_{k=1}^{n} \frac{1}{\hat{\theta}_2} + \sum_{k=1}^{n} \frac{(x_k - \hat{\theta}_1)^2}{\hat{\theta}_2^2} = 0 \qquad (2) \implies \hat{\theta}_2 = \hat{\sigma}^2 = \frac{1}{n} \sum_{k=1}^{n} (x_k - \hat{\mu})^2$$