

Chapter 3:

Maximum-Likelihood and Bayesian Parameter Estimation

Introduction

- Data availability in a Bayesian framework
 - We could design an optimal classifier if we knew:
 - $P(\omega_i)$ (priors)
 - $P(x \mid \omega_i)$ (class-conditional densities)

Unfortunately, we rarely have this complete information

- Design a classifier from a training sample
 - No problem with prior estimation
 - Samples are often too small for class-conditional estimation (large dimension of feature space!)

- Apriori information about the problem
 - Normality of $P(x \mid \omega_i)$

$$P(x \mid \omega_i) \sim N(\mu_i, \Sigma_i)$$

- Characterized by 2 parameters

- Estimation techniques
 - Maximum-Likelihood (ML) and Bayesian estimation
 - Results are nearly identical, but the approaches are different
- Parameters in ML estimation are fixed but unknown
- Bayesian methods view the parameters as random variables having some known distribution

- Best parameters are obtained by maximizing the probability of obtaining the samples observed
- In either approach, we use $P(\omega_i | x)$ for our classification rule!

Maximum-Likelihood Estimation

- Has good convergence properties as the sample size increases
- Simpler than any other alternative techniques
- General principle
 - Assume we have c classes and
 - $P(x | \omega_j) \sim N(\mu_j, \Sigma_j)$
 - $P(x | \omega_j) \equiv P(x | \omega_j, \theta_j)$ where:

$$\theta_j = (\mu_j, \Sigma_j) = (\mu_j^1, \mu_j^2, \dots, \sigma_j^{11}, \sigma_j^{22}, \text{cov}(x_j^m, x_j^n) \dots)$$

Maximum-Likelihood Estimation

- Use the information provided by the training samples to estimate $\theta = (\theta_1, \theta_2, \dots, \theta_c)$, each θ_i ($i = 1, 2, \dots, c$) is associated with each category
- Suppose that D contains n samples, x_1, x_2, \dots, x_n

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

$P(D | \theta)$ is called the likelihood of θ w.r.t. the set of samples)

- ML estimate of θ is, by definition the value that maximizes $P(D | \theta)$
- “It is the value of θ that best agrees with the actually observed training sample”

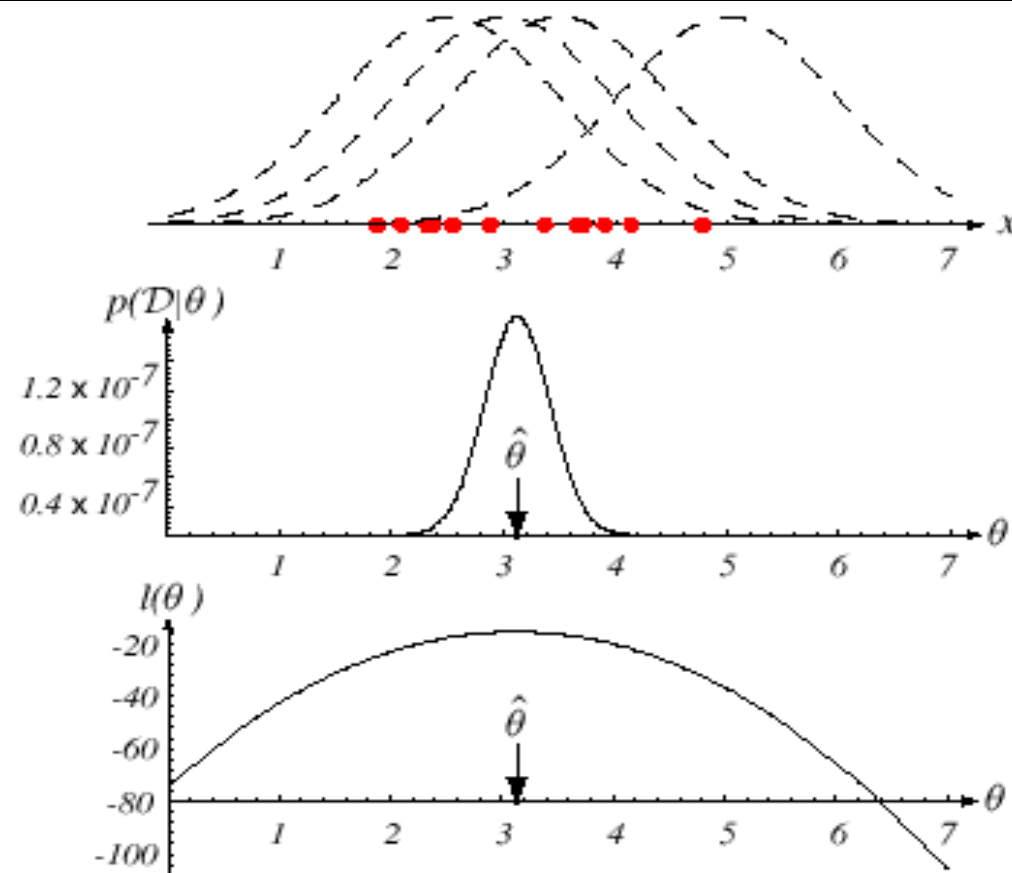


FIGURE 3.1. The top graph shows several training points in one dimension, known or assumed to be drawn from a Gaussian of a particular variance, but unknown mean. Four of the infinite number of candidate source distributions are shown in dashed lines. The middle figure shows the likelihood $p(\mathcal{D}|\theta)$ as a function of the mean. If we had a very large number of training points, this likelihood would be very narrow. The value that maximizes the likelihood is marked $\hat{\theta}$; it also maximizes the logarithm of the likelihood—that is, the log-likelihood $l(\theta)$, shown at the bottom. Note that even though they look similar, the likelihood $p(\mathcal{D}|\theta)$ is shown as a function of θ whereas the conditional density $p(x|\theta)$ is shown as a function of x . Furthermore, as a function of θ , the likelihood $p(\mathcal{D}|\theta)$ is not a probability density function and its area has no significance. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

- Optimal estimation

- Let $\theta = (\theta_1, \theta_2, \dots, \theta_p)^t$ and let ∇_{θ} be the gradient operator

$$\nabla_{\theta} = \left[\frac{\partial}{\partial \theta_1}, \frac{\partial}{\partial \theta_2}, \dots, \frac{\partial}{\partial \theta_p} \right]^t$$

- We define $l(\theta)$ as the log-likelihood function $l(\theta) = \ln P(D | \theta)$
 - New problem statement: “Determine θ that maximizes the log-likelihood”

$$\hat{\theta} = \arg \max_{\theta} l(\theta)$$

- Set of necessary conditions for an optimum is:

$$\nabla_{\theta} l = \sum_{k=1}^{k=n} \nabla_{\theta} \ln P(x_k | \theta)$$
$$\nabla_{\theta} l = 0$$

Example of a specific case: Unknown μ

- $P(x_i | \mu) \sim N(\mu, \Sigma)$
- (Samples are drawn from a multivariate normal population)

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

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

- $\theta = \mu$ therefore, the ML estimate for μ must satisfy:

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

- Multiplying by Σ and rearranging, we obtain:

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

- Just the arithmetic average of the samples of the training data
- Conclusion:
 - If $P(x_k | \omega_j)$ ($j = 1, 2, \dots, c$) is supposed to be Gaussian in a d -dimensional feature space; then we can estimate the vector $\theta = (\theta_1, \theta_2, \dots, \theta_c)^t$ and perform an optimal classification!

Gaussian Case: Unknown μ and σ

- $\theta = (\theta_1, \theta_2) = (\mu, \sigma^2)$

$$l = \ln P(x_k | \theta) = -\frac{1}{2} \ln 2\pi\theta_2 - \frac{1}{2\theta_2} (x_k - \theta_1)^2$$
$$\nabla_{\theta} l = \begin{pmatrix} \frac{\partial}{\partial \theta_1} (\ln P(x_k | \theta)) \\ \frac{\partial}{\partial \theta_2} (\ln P(x_k | \theta)) \end{pmatrix}$$
$$= \begin{cases} \frac{1}{\theta_2} (x_k - \theta_1) \\ -\frac{1}{2\theta_2} + \frac{(x_k - \theta_1)^2}{2\theta_2^2} \end{cases}$$

- Summation:

$$\left\{ \begin{array}{l} \sum_{k=1}^{k=n} \frac{1}{\hat{\theta}_2} (x_k - \theta_1) = 0 \quad (1) \\ -\sum_{k=1}^{k=n} \frac{1}{\hat{\theta}_2} + \sum_{k=1}^{k=n} \frac{(x_k - \hat{\theta}_1)^2}{\hat{\theta}_2^2} = 0 \quad (2) \end{array} \right.$$

- Combining (1) and (2):

$$\mu = \sum_{k=1}^{k=n} \frac{x_k}{n} \quad ; \quad \sigma^2 = \frac{\sum_{k=1}^{k=n} (x_k - \mu)^2}{n}$$

Bias

- ML estimate for σ^2 is biased, i.e., Expected value of sample variance over all datasets of size n is not equal to true variance

$$E\left[\frac{1}{n}\sum(x_i - \bar{x})^2\right] = \frac{n-1}{n}\sigma^2 \neq \sigma^2$$

- An elementary unbiased estimator for Σ is:

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

Sample covariance matrix

- If an estimator is unbiased for all distributions, then it is called *absolutely unbiased*
- If the estimator tends to become unbiased as the number of samples become very large, then the estimator is *asymptotically unbiased*

- Bayesian Estimation (BE)
- Bayesian Parameter Estimation: General Estimation
- Bayesian Parameter Estimation: Gaussian Case

Bayesian Estimation

- Bayesian learning to pattern classification problems
 - In MLE, θ was supposed fix
 - In BE, θ is a random variable
 - The computation of posterior probabilities $P(\omega_i | x)$ lies at the heart of Bayesian classification
 - Goal: Compute $P(\omega_i | x, D)$
 - Given the sample D , Bayes formula can be written

$$P(\omega_i | x, D) = \frac{P(x | \omega_i, D) \cdot P(\omega_i | D)}{\sum_{j=1}^c P(x | \omega_j, D) \cdot P(\omega_j | D)}$$

Bayesian Estimation ...

- To demonstrate the preceding equation, use:

$$P(x, D | \omega_i) = P(x | D, \omega_i) \cdot P(D | \omega_i)$$

$$P(x | D) = \sum_j P(x, \omega_j | D)$$

$$P(\omega_i) = P(\omega_i | D) \quad (\text{Training sample provides this!})$$

Thus:

$$P(\omega_i | x, D) = \frac{P(x | \omega_i, D) \cdot P(\omega_i)}{\sum_{j=1}^c P(x | \omega_j, D) \cdot P(\omega_j)}$$

Bayesian Parameter Estimation: Gaussian Case

- Goal: Estimate θ using the a-posteriori density $P(\theta | D)$
 - The univariate case: $P(\mu | D)$
 - μ is the only unknown parameter

$$p(x | \mu) \sim N(\mu, \sigma^2)$$
$$p(\mu) \sim N(\mu_0, \sigma_0^2)$$

- (μ_0 and σ_0 are known!)

$$P(\mu | \mathbf{D}) = \frac{P(\mathbf{D} | \mu) \cdot P(\mu)}{\int P(\mathbf{D} | \mu) \cdot P(\mu) d\mu}$$

$$= \alpha \prod_{k=1}^{k=n} P(x_k | \mu) \cdot P(\mu)$$

$$P(\mu | \mathbf{D}) \sim N(\mu_n, \sigma_n^2)$$

$$P(\mu) \sim N(\mu_0, \sigma_0^2)$$

$$p(\mu | \mathcal{D}) = \alpha \prod_{k=1}^n \overbrace{\frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{1}{2} \left(\frac{x_k - \mu}{\sigma} \right)^2 \right]}^{p(x_k | \mu)} \overbrace{\frac{1}{\sqrt{2\pi}\sigma_0} \exp \left[-\frac{1}{2} \left(\frac{\mu - \mu_0}{\sigma_0} \right)^2 \right]}^{p(\mu)}$$

$$= \alpha' \exp \left[-\frac{1}{2} \left(\sum_{k=1}^n \left(\frac{\mu - x_k}{\sigma} \right)^2 + \left(\frac{\mu - \mu_0}{\sigma_0} \right)^2 \right) \right]$$

$$= \alpha'' \exp \left[-\frac{1}{2} \left[\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right) \mu^2 - 2 \left(\frac{1}{\sigma^2} \sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2} \right) \mu \right] \right],$$

$$\begin{aligned}
 p(\mu|\mathcal{D}) &= \alpha \prod_{k=1}^n \overbrace{\frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{1}{2} \left(\frac{x_k - \mu}{\sigma} \right)^2 \right]}^{p(x_k|\mu)} \overbrace{\frac{1}{\sqrt{2\pi}\sigma_0} \exp \left[-\frac{1}{2} \left(\frac{\mu - \mu_0}{\sigma_0} \right)^2 \right]}^{p(\mu)} \\
 &= \alpha' \exp \left[-\frac{1}{2} \left(\sum_{k=1}^n \left(\frac{\mu - x_k}{\sigma} \right)^2 + \left(\frac{\mu - \mu_0}{\sigma_0} \right)^2 \right) \right] \\
 &= \alpha'' \exp \left[-\frac{1}{2} \left[\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right) \mu^2 - 2 \left(\frac{1}{\sigma^2} \sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2} \right) \mu \right] \right],
 \end{aligned}$$

$$p(\mu|\mathcal{D}) = \frac{1}{\sqrt{2\pi}\sigma_n} \exp \left[-\frac{1}{2} \left(\frac{\mu - \mu_n}{\sigma_n} \right)^2 \right]$$

$$\frac{\mu_n}{\sigma_n^2} = \frac{n}{\sigma^2} \bar{x}_n + \frac{\mu_0}{\sigma_0^2},$$

$$\frac{1}{\sigma_n^2} = \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}$$

$$\bar{x}_n = \frac{1}{n} \sum_{k=1}^n x_k$$

$$\mu_n = \left(\frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \bar{x}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0$$

$$\sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2}.$$

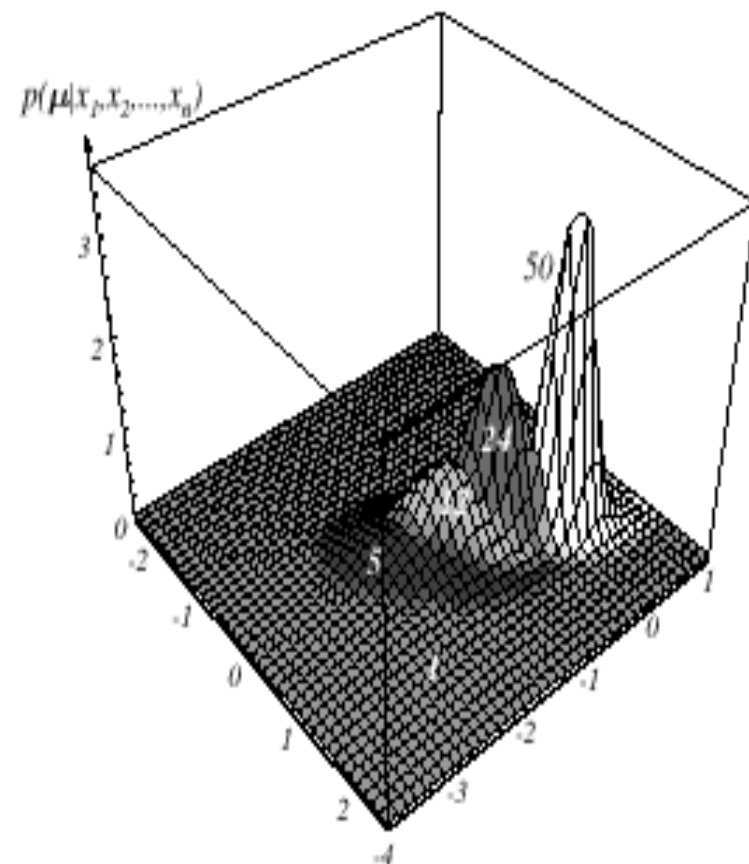
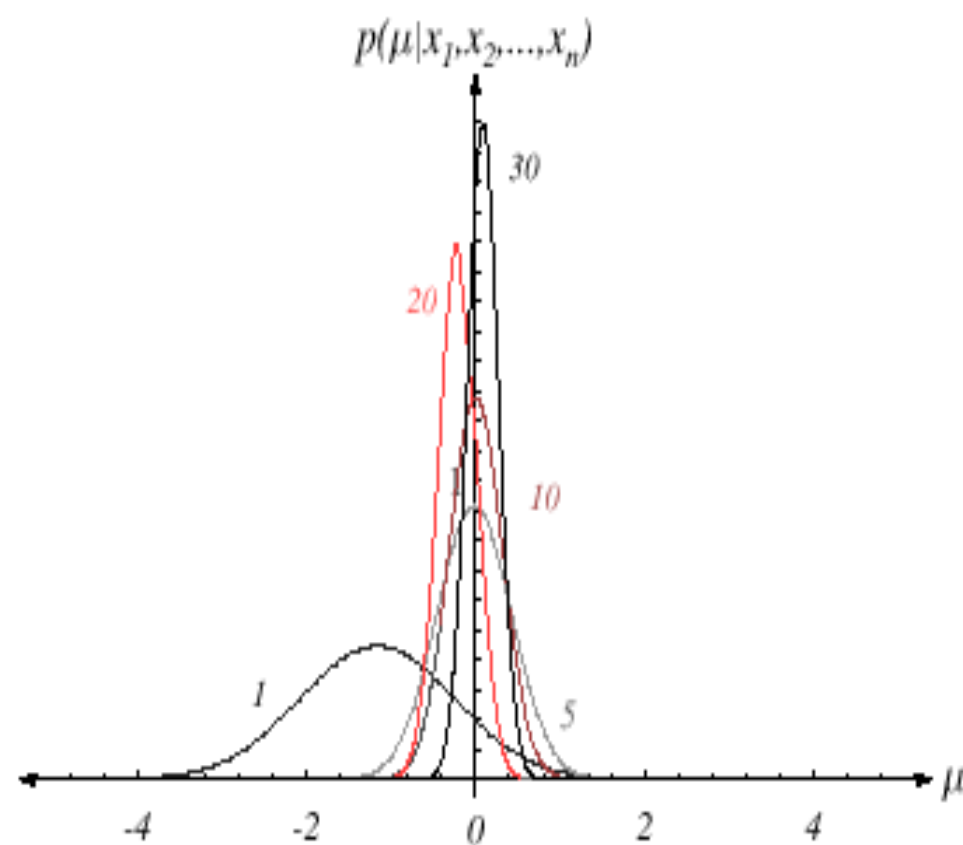


FIGURE 3.2. Bayesian learning of the mean of normal distributions in one and two dimensions. The posterior distribution estimates are labeled by the number of training samples used in the estimation. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Univariate Case

– The univariate case $P(x \mid D)$

- $P(\mu \mid D)$ computed
- $P(x \mid D)$ remains to be computed!

$$P(x \mid D) = \int P(x \mid \mu) P(\mu \mid D) d\mu \text{ is Gaussian}$$

- It provides: $P(x \mid D) \sim N(\mu_n, \sigma^2 + \sigma_n^2)$
- (Desired class-conditional density $P(x \mid D_j, \omega_j)$)
- Therefore: $P(x \mid D_j, \omega_j)$ together with $P(\omega_j)$
- And using Bayes formula for classification

Bayesian Parameter Estimation: General Theory

- $P(x \mid D)$ computation can be applied to any situation in which the unknown density can be parameterized: the basic assumptions are:
 - The form of $P(x \mid \theta)$ is assumed known, but the value of θ is not known exactly
 - Our knowledge about θ is assumed to be contained in a known prior density $P(\theta)$
 - The rest of our knowledge θ is contained in a set D of n random variables x_1, x_2, \dots, x_n that follows $P(x)$

Bayesian Parameter Estimation: General Theory...

- The basic problem is: “Compute the posterior density $P(\theta \mid D)$ ”
- then “Derive $P(x \mid D)$ ”

- Using Bayes formula:
$$P(\theta \mid D) = \frac{P(D \mid \theta).P(\theta)}{\int P(D \mid \theta).P(\theta)d\theta},$$

- And by independence assumption:

$$P(D \mid \theta) = \prod_{k=1}^{k=n} P(x_k \mid \theta)$$

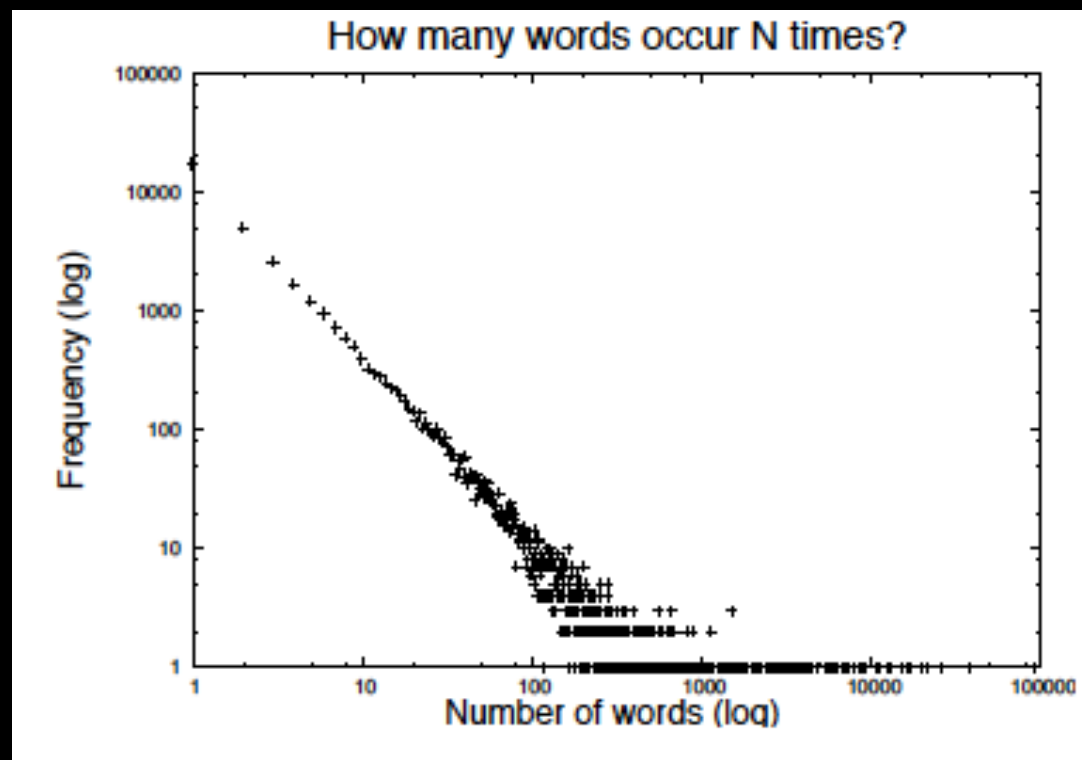
Three Sources of Error in Classification

- Bayes or Indistinguishability Error
 - Due to overlapping densities. Inherent property of given feature set
- Model Error
 - Due to an incorrect model. Can only be eliminated if designer specifies true model that generated the data.
- Estimation Error
 - Parameters are estimated from a finite sample.

Smoothing

How much probability does the model assign to strings/ pixel values it hasn't seen before?

- An empirical fact about language:
- A small number of events occur with high frequency
- A large number of events occur with low frequency



Zipf's Law

Smoothing

- Add-1 (Laplace) smoothing: Assume every seen or unseen event occurred once more than it did in the training data

$$\begin{array}{ll} \text{MLE} & P(w_i) = \frac{C(w_i)}{\sum_j C(w_j)} = \frac{C(w_i)}{N} \\ \text{Add One} & P(w_i) = \frac{C(w_i)+1}{\sum_j (C(w_j)+1)} = \frac{C(w_i)+1}{N+V} \end{array}$$

- Disadvantage of Add-1 Smoothing
 - Takes away too much probability mass from seen events.
 - Assigns too much total probability mass to unseen events.

Estimators

$$P(X) = \frac{N_X + f}{T + E_\Omega \cdot f},$$

- f is a fudge factor
- Expected Likelihood Estimator sets $f = 0.5$
- Laplace and add-one estimators are the same thing and set $f = 1$,
- add-tiny sets $f = 1/T$

Good Turing Smoothing

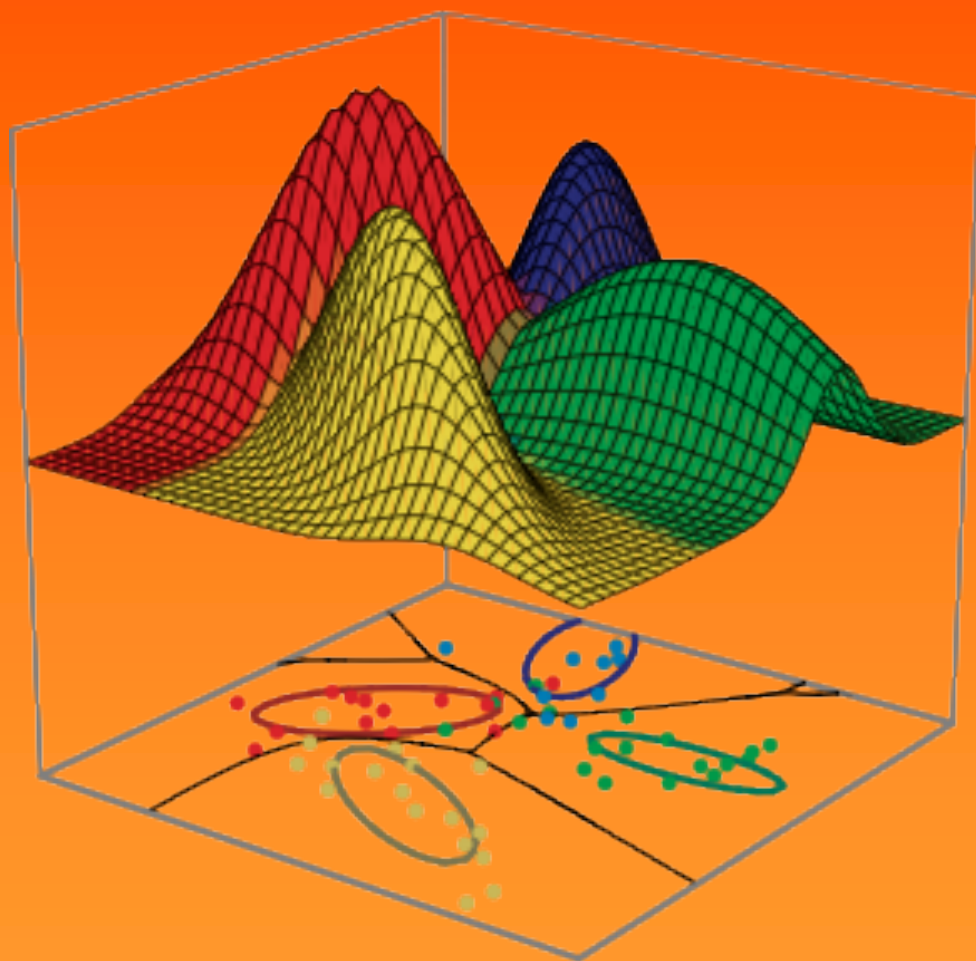
$$F_X = \frac{(N_X + 1)}{T} \cdot \frac{E(N_X + 1)}{E(N_X)}$$

- X is the event, N_X is the number of times you have seen event X , T is the sample size and $E(n)$ is an estimate of how many different events happened exactly n times. Translating that into text-analysis terms, X is a word, N_X is the number of times you have seen word X , T is the size of the corpus and $E(n)$ is an estimate of how many different words were observed exactly n times **3**.

Good Turing Smoothing

- Take a corpus of 30000 English words, our universe is all English words, our event, X , is the word “unusualness.” The word “unusualness” appears once, so $NX = 1$. In a reasonable corpus, you might have 10000 different words that appear once, so $E(1) = 10000$, and you might have 3000 words that appear twice, giving $E(2) = 3000$. The Good-Turing estimate of the probability of “unusualness” is then

$$P(\text{unusualness}) = \frac{2}{30000} \cdot \frac{3000}{10000},$$



Pattern Classification

All materials in these slides were taken from
Pattern Classification (2nd ed) by R. O. Duda, P. E. Hart and D. G. Stork, John Wiley & Sons, 2000
with the permission of the authors and the publisher