

(1)
Alpaydin
chp 4.

Parametric Methods

3/26/2008

How to estimate the probabilities to make decisions (previous lecture).

Abstract: estimate the sufficient statistics which are functions of the data samples. These are sufficient to determine the parameters of the distribution.

Density estimator $p(x)$
Use this to estimate $p(x|\theta)$ & $p(\theta)$

Maximum Likelihood Estimation (MLE)

Sample $X = \{x^t\}_{t=1}^N$

i.i.d. from some distribution $p(x|\theta)$

likelihood $l(\theta) \equiv p(X|\theta) = \prod_{t=1}^N p(x^t|\theta)$ parameter.

MLE: Find $\hat{\theta}$ to maximize $l(\theta)$

or, equivalently maximize the log likelihood

$$L(\theta|X) \equiv \log l(\theta|X) = \sum_{t=1}^N \log p(x^t|\theta)$$

(2)

Bernoulli Density

two outcomes,
random variable X

$X = 1$, with prob p
 $X = 0$, with prob $1-p$.

$$P(X) = p^x (1-p)^{1-x}, \quad x \in \{0, 1\}.$$

p is the only parameter.

Log-likelihood

$$L(p|X) = \log \prod_{t=1}^N p^{x^t} (1-p)^{1-x^t}$$

$$= \sum_t x^t \log p + (N - \sum_t x^t) \log(1-p).$$

Solve $\frac{\partial L(p|X)}{\partial p} = 0$

gives $\hat{p} = \frac{\sum_t x^t}{N}$.

Intuition, no. of positive outcomes
divided by number of samples.

Multinomial Density

generalization of Bernoulli

K classes, probability p_i , $\sum_{i=1}^K p_i = 1$

$$P(x_1, x_2, \dots, x_K) = \prod_{i=1}^K p_i^{x_i}$$

x_i 's indicator variables

$x_i = 1$, if outcome is i
 0 , otherwise.

(3)

N experiments with outcomes $X = \{x^t\}_{t=1}^N$
 $x^t_i = \begin{cases} 1, & \text{if experiment } t \text{ chose state } i \\ 0, & \text{otherwise} \end{cases}$

MLE of P_i is $\hat{P}_i = \frac{\sum_t x^t_i}{N}$
proportion of sample which lie in class i .

Gaussian (Normal) Density

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad -\infty < x < \infty$$

$\mathcal{N}(\mu, \sigma^2)$

Sample $X = \{x^t\}_{t=1}^N$, with $x^t \sim \mathcal{N}(\mu, \sigma^2)$

log likelihood $\mathcal{L}(\mu, \sigma | X) = -\frac{N}{2} \log 2\pi - N \log \sigma - \sum_t \frac{(x^t - \mu)^2}{2\sigma^2}$

MLE $m = \frac{\sum_t x^t}{N}$
 $s^2 = \frac{\sum_t (x^t - m)^2}{N}$

Convention: Greek for the parameters.
Roman for their estimates.

(page 4)

Exponential Distributions and Maximum Likelihood

$$P(x|\lambda) = \frac{e^{\lambda \cdot \phi(x)}}{Z[\lambda]}$$

exponential distribution
statistics $\phi(\cdot)$

parameter λ

$$Z[\lambda] = \sum_x e^{\lambda \cdot \phi(x)}, \text{ if } x \text{ is discrete-valued}$$

$$= \int dx e^{\lambda \cdot \phi(x)}, \text{ if } x \text{ is continuous-valued}$$

normalization
term.

Ensures $\sum_x P(x|\lambda) = 1$, or $\int dx P(x|\lambda) = 1$.

Note: (1) most "named" distributions can be expressed as exponential distributions \rightarrow e.g. Gaussian, Poisson, Bernoulli, ...

(2) any distribution can be approximated arbitrarily accurately by an exponential distribution.

(3) mixture models can be expressed as exponential distributions by introducing hidden variables

Example: Gaussian in 1-D $P(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$

$$\log P(x|\mu, \sigma) = \frac{-1}{2\sigma^2} x^2 + \frac{x\mu}{\sigma^2} - \frac{1}{2\sigma^2} \mu^2 - \log \sqrt{2\pi}\sigma$$

$$\log P(x|\lambda) = \lambda \cdot \phi(x) - \log Z[\lambda] \quad \text{set } \phi(x) = (x, x^2)$$

$$\lambda = \left(-\frac{1}{2\sigma^2}, \frac{\mu}{\sigma^2} \right)$$

$$\log Z[\lambda] = -\frac{1}{2\sigma^2} \mu^2 - \log \sqrt{2\pi}\sigma$$

(5)

Note: the sufficient statistics $\phi(x)$ are the properties of the data x that we need to store (we can throw out the rest of the data).

For a set of samples x^1, x^2, \dots, x^M from a Gaussian,

we only need to store the statistics

$$\frac{1}{M} \sum_{i=1}^M (x^i) \quad \text{and} \quad \frac{1}{M} \sum_{i=1}^M (x^i)^2.$$

the parameters of the Gaussian — μ & σ^2 — estimated by ML will depend only on these statistics:

e.g. we don't need to also store

$$\frac{1}{M} \sum_{i=1}^M (x^i)^4, \quad \text{or even all the samples } x^1, x^2, \dots, x^M)$$

Another example: Bernoulli.

$$p(x) = \theta^x (1-\theta)^{1-x}$$

$$\log p(x) = x \log \theta + (1-x) \log (1-\theta)$$

$$= x \{ \log \theta - \log (1-\theta) \} + \log (1-\theta)$$

Exponential

$$\phi(x) = x$$

$$\lambda = \log \theta - \log (1-\theta)$$

$$\log \zeta(\lambda) = \log (1-\theta)$$

(6) Maximum Likelihood for Exponential Distribution

Samples $\{x^i : i=1 \dots M\}$ from $P(x|\lambda) = \frac{e^{-\lambda \phi(x)}}{Z(\lambda)}$

ML say maximize $\prod_{i=1}^M P(x^i|\lambda)$
w.r.t. λ

$$\hat{\lambda} = \text{ARG MAX}_{\lambda} \sum_{i=1}^M \log P(x^i|\lambda)$$

$$\hat{\lambda} = \text{ARG MAX}_{\lambda} \left(\lambda \sum_{i=1}^M \phi(x^i) - M \log Z(\lambda) \right)$$

Equivalently: minimize $F(\lambda) = \log Z(\lambda) - \lambda \frac{1}{M} \sum_{i=1}^M \phi(x^i)$

where $\psi = \frac{1}{M} \sum_{i=1}^M \phi(x^i)$ is the empirical statistic
e.g. $\left(\frac{1}{M} \sum_{i=1}^M x^i, \frac{1}{M} \sum_{i=1}^M (x^i)^2 \right)$ for Gaussian

Minimize: solve $\frac{\partial F(\lambda)}{\partial \lambda} = 0$, for λ :

Note $F(\lambda)$ is convex - $\frac{\partial^2 F}{\partial \lambda^2} \geq 0$ by properties of $\log Z(\lambda)$

so there is a unique solution

(property of $\log Z(\lambda)$) $\frac{\partial F}{\partial \lambda} = \frac{\partial \log Z(\lambda)}{\partial \lambda} - \psi = \sum_x P(x|\lambda) \phi(x) - \psi$

Maximum likelihood selects λ for that the expected statistic $\sum_x P(x|\lambda) \phi(x)$ equal the observed statistic ψ

- e.g. for a Gaussian the mean and variance parameters λ, μ are chosen to equal the mean and variance of the data

(7)

Important:

how to solve $\sum_x P(x|\lambda) \underline{\phi}(x) = \underline{\psi}$?
 (or $\int dx P(x|\lambda) \underline{\phi}(x) = \underline{\psi}$) ?

For special distributions — like Gaussians, Bernoulli, Poisson — these can be solved analytically to give a closed form solution

→ e.g. $\mu = \frac{1}{N} \sum_{i=1}^N x^i$ for Gaussians
 $\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x^i)^2 + \mu^2$

But for other distributions we cannot solve analytically.
 (or Generalized Iterative Scaling GIS)

In these cases we can do steepest descent on $F(\lambda)$ — which is convex, so steepest descent converges to the unique solution.

$$\underline{\lambda}^{t+1} = \underline{\lambda}^t - \epsilon \left. \frac{\partial F(\lambda)}{\partial \lambda} \right|_{\underline{\lambda}^t} \quad \begin{matrix} t=0 \\ t - \text{index of} \\ \text{iteration.} \end{matrix}$$

$$\underline{\lambda}^{t+1} = \underline{\lambda}^t - \epsilon \sum_x \underline{\phi}(x) P(x; \underline{\lambda}^t) + \epsilon \underline{\psi}$$

converges — when $\sum_x \underline{\phi}(x) P(x; \underline{\lambda}^t) = \underline{\psi}$

Problem: each iteration

requires estimating $\sum_x \underline{\phi}(x) P(x; \underline{\lambda}^t)$

model statistics
 balance the empirical statistics

which is difficult — note, if we can compute this analytically then we can probably solve $\sum_x \underline{\phi}(x) P(x; \underline{\lambda}) = \underline{\psi}$ directly (e.g. for Gaussian $\int \underline{\phi}(x) P(x; \lambda) dx = (\mu, \sigma^2 = \mu^2)$)

(page 8)

Alternate perspective: Maximum Entropy.

Entropy of a distribution

$$-\sum_x p(x) \log p(x), \text{ or } -\int dx p(x) \log p(x)$$

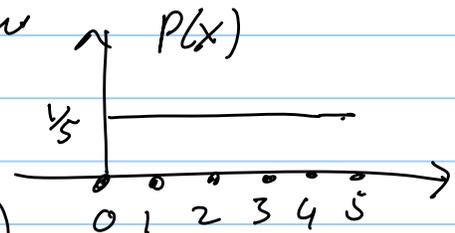
expected information from a sample x of a distribution

Shannon's Theory of Information

uniform distribution
entropy

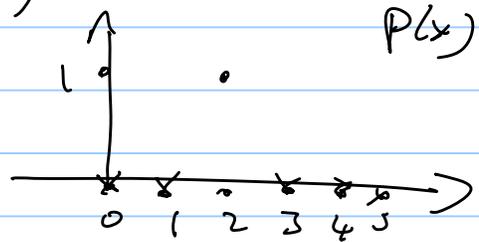
$$-5 \left(\frac{1}{5}\right) \log\left(\frac{1}{5}\right) = \log 5$$

generally $\log N$ (if N no. of states)



entropy

$$-4(0 \log 0) - (1 \log 1) = 0$$



$$\left(\begin{array}{l} 0 \log 0 = 0 \\ 1 \log 1 = 0 \end{array} \right)$$

no information is transmitted if we get a sample from this distribution, since we know what the value will be before we see it (it has to be 2)

But we do get information from a sample from the uniform distribution

Shannon; says we should encode a sample x from $p(x)$ by $-\log p(x)$ bits — so samples with high probability have low number of bits — then the expected coding cost (in bits) is just the entropy

$$\sum_x p(x) (-\log p(x)) = -\sum_x p(x) \log p(x)$$

(page 9)

Maximum entropy principle:

suppose that we observe some statistics $\phi(x)$ and they take a value $\underline{\psi}$

(e.g. the mean of the data takes a value).

What distribution describes the data?

Maximum Entropy Principle (see Jaynes)

Suppose we select the distribution $P(x)$ that maximizes the entropy (i.e. is as uniform as possible) but is consistent with the data

$$L(P; \lambda, \mu) = - \sum_x P(x) \log P(x) + \mu \left(\sum_x P(x) - 1 \right) + \lambda \left(\sum_x P(x) \phi(x) - \psi \right)$$

Lagrange multipliers
imposed on the constraint

$$\frac{\partial L}{\partial P} = 0 \Rightarrow -\log P(x) + 1 + \mu + \lambda \cdot \phi(x) = 0$$
$$\Rightarrow P(x) = e^{\lambda \cdot \phi(x)} e^{1+\mu} \quad (*)$$

$$\frac{\partial L}{\partial \mu} = 0 \Rightarrow \sum_x P(x) = 1, \text{ normalized}$$

$$\frac{\partial L}{\partial \lambda} = 0 \Rightarrow \sum_x P(x) \phi(x) = \psi, \text{ consistent with the data}$$

Imposing normalization on (*) gives

$$P(x; \lambda) = \frac{1}{Z(\lambda)} e^{\lambda \cdot \phi(x)} \rightarrow \text{exponential distribution!}$$

$$\text{s.t. } \sum_x P(x; \lambda) \phi(x) = \psi$$

So Maximum Likelihood with Exponential Distributions is the same as doing Maximum Entropy with the statistics

(10) Maximum Entropy suggests a solution to the problem - how do we know what distribution to use?

Maximum Entropy strategy is to search over the statistics - instead of searching over distributions

i.e. statistic $\phi_1(\cdot)$ gives distribution $P(x) = \frac{1}{Z(\lambda_1)} e^{\lambda_1 \phi_1(x)}$

" " $\phi_2(\cdot)$ " " $P(x) = \frac{1}{Z(\lambda_2)} e^{\lambda_2 \phi_2(x)}$

" $\phi_1(x) \& \phi_2(x)$ " " $P(x) = \frac{1}{Z(\lambda_1, \lambda_2)} e^{\lambda_1 \phi_1(x) + \lambda_2 \phi_2(x)}$

"Searching over Statistics" is complicated and difficult
- see papers by Della Pietra, Della Pietra, Elftadey
Zhu, Wu, Mumford.

Beyond scope of this course.