

# Lecture 4

Spring 2014

Note Title

4/9/2014

## TOPICS:

- (I) Learning probability distributions by Maximum Likelihood (ML)
- (II) Exponential distributions, Sufficient Statistics, and ML learning.
- (III) Kullback-Leibler divergence, learning "approximate" distributions.
- (IV) Advanced Topics,  
Maximum Entropy Principle.  
Model Pursuit.  
Model Selection.

(2) Maximum Likelihood (ML) learning of parametric probability distributions. (non-parametric, later lecture)

Note Title

Learn distributions like  $P(x|y)$  likelihood function or  $P(y)$  prior.

For simplicity, we just study learning  $P(x)$ .  
(Note: learning  $P(y|x)$  - regression - later in the course.)

Parametric Distribution  $\rightarrow P(x|\theta)$  parameter

E.g. Gaussian  $P(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ ,  $\theta = (\mu, \sigma)$

Data  $\mathcal{X}_N = \{x_i : i=1 \dots N\}$

Assume data is i.i.d. from unknown  $P(x)$

Maximum Likelihood (ML):

$$P(\mathcal{X}_N|\theta) = \prod_{i=1}^N P(x_i|\theta) \quad (\text{i.i.d. assumption})$$

Set  $\hat{\theta} = \underset{\theta}{\text{ARG MAX}} P(\mathcal{X}_N|\theta)$   
 $P(\mathcal{X}_N|\hat{\theta}) \geq P(\mathcal{X}_N|\theta)$  for any  $\theta$

This is equivalent to  $\underset{\theta}{\text{arg max}} \log P(\mathcal{X}_N|\theta)$

$$\hat{\theta} = \underset{\theta}{\text{ARG MAX}} \left\{ \sum_{i=1}^N \log P(x_i|\theta) \right\}$$

For Gaussian:  $\sum_{i=1}^N \log P(x_i|\mu, \sigma) = -\sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} - \sum_{i=1}^N \log \sqrt{2\pi}\sigma$

To estimate  $(\hat{\mu}, \hat{\sigma})$  differentiate w.r.t.  $\mu, \sigma$  (ie. maximize  $\log P(\mathcal{X}_N|\mu, \sigma)$  w.r.t.  $\mu, \sigma$ .)

$$\frac{\partial}{\partial \mu} \log P(\mathcal{X}_N|\mu, \sigma) = \frac{1}{\sigma^2} \sum_{i=1}^N (x_i - \mu)$$

$$\frac{\partial}{\partial \sigma} \log P(\mathcal{X}_N|\mu, \sigma) = \frac{1}{\sigma^3} \sum_{i=1}^N (x_i - \mu)^2 - N/\sigma$$

Hence solution  $\left. \begin{aligned} \hat{\mu} &= \frac{1}{N} \sum_{i=1}^N x_i \\ \hat{\sigma}^2 &= \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2 \end{aligned} \right\}$

(3) Note: the Gaussian is a special case because it gives a simple analytic formula for  $\hat{\mu}, \hat{\sigma}^2$ .

This illustrates the ML estimator.

Note: What about MAP?

Why not use a prior  $P(\theta)$ ?

$$\begin{aligned}\hat{\theta}_{\text{MAP}} &= \underset{\theta}{\text{ARG MAX}} P(\mathcal{X}_N | \theta) P(\theta) \\ &= \underset{\theta}{\text{ARG MAX}} \left\{ \overbrace{\log P(\mathcal{X}_N)}^{\text{width of } \theta} + \log P(\theta) \right\} \\ &= \underset{\theta}{\text{ARG MAX}} \left\{ \underbrace{\sum_{i=1}^N \log P(x_i | \theta)}_{N \text{ data terms}} + \underbrace{\log P(\theta)}_{1 \text{ prior term}} \right\}\end{aligned}$$

If  $N$  is large, then the prior will have little effect (except for special cases).

Note: we can also use loss functions and all the machinery of Bayes Decision Theory.

In practice, loss functions are often not used when learning probability distributions - but they are used for learning classifiers (later in course)

## (4) Topic II: Exponential Distributions

### Exponential Distributions

$$P(\underline{x}|\underline{\lambda}) = \frac{1}{Z[\underline{\lambda}]} e^{\underline{\lambda} \cdot \underline{\phi}(\underline{x})}$$

normalization factor.

$\underline{\lambda}$  - parameters  
 $\underline{\phi}(\underline{x})$  - statistics.

$\underline{\lambda} = (\lambda_1, \lambda_2, \dots, \lambda_M)$   
 $\underline{\phi}(\underline{x}) = (\phi_1(\underline{x}), \phi_2(\underline{x}), \dots, \phi_M(\underline{x}))$

Almost every named distribution can be expressed as an exponential distribution.

For Gaussian in 1-dimension.

write  $\underline{\phi}(\underline{x}) = (x, x^2)$      $\underline{\lambda} = \lambda_1, \lambda_2$

$$P(x|\underline{\lambda}) = \frac{1}{Z[\underline{\lambda}]} e^{\lambda_1 x + \lambda_2 x^2}$$

compare to  $\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

Translation

$$\begin{cases} \lambda_2 = -1/2\sigma^2 \\ \lambda_1 = \mu/\sigma^2 \\ Z[\underline{\lambda}] = \sqrt{2\pi}\sigma e^{\mu^2/2\sigma^2} \end{cases}$$

Similar translation into exponential distributions can be made for Poisson, Beta, Dirichlet ~ most (all) distributions you have been taught.

(5)

## Learning an Exponential Distribution

You can learn them by Maximum Likelihood,

Examples:

$$P(\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N | \underline{\lambda}) = \prod_{n=1}^N \frac{e^{\underline{\lambda} \cdot \underline{\phi}(\underline{x}_n)}}{Z(\underline{\lambda})}$$

Maximize w.r.t.  $\underline{\lambda}$  ||

This has a very nice form, which occurs because the exponential distribution depends on the data  $\underline{x}$  only in terms of the function  $\underline{\phi}(\underline{x})$  — the sufficient statistics.

Note: Important factor. The normalization term  $Z(\underline{\lambda})$  is a function of  $\underline{\lambda}$ ,  $Z(\underline{\lambda}) = \sum_{\underline{x}} e^{\underline{\lambda} \cdot \underline{\phi}(\underline{x})}$

Claim:  $\frac{\partial \log Z(\underline{\lambda})}{\partial \underline{\lambda}} = \sum_{\underline{x}} \underline{\phi}(\underline{x}) P(\underline{x} | \underline{\lambda})$   
expected value of the statistics  $\underline{\phi}(\underline{x})$  w.r.t.  $P(\underline{x} | \underline{\lambda})$

Proof:  $\frac{\partial \log Z(\underline{\lambda})}{\partial \underline{\lambda}} = \frac{1}{Z(\underline{\lambda})} \frac{\partial Z(\underline{\lambda})}{\partial \underline{\lambda}} = \frac{1}{Z(\underline{\lambda})} \sum_{\underline{x}} \underline{\phi}(\underline{x}) e^{\underline{\lambda} \cdot \underline{\phi}(\underline{x})} = \sum_{\underline{x}} \underline{\phi}(\underline{x}) P(\underline{x} | \underline{\lambda})$ .

(6)

Claim: For exponential distribution, ML corresponds to finding the value of  $\underline{\lambda}$  st. the model statistics are equal to the data statistics  
solve  $\sum_{\underline{x}} \underline{\phi}(\underline{x}) P(\underline{x} | \underline{\lambda}) = \frac{1}{N} \sum_{i=1}^N \underline{\phi}(x_i)$ .

E.G. For Gaussian.  
model statistics  $\int d\underline{x} \underline{x} \frac{1}{(2\pi)^{N/2} |\det \underline{\Sigma}|} e^{-\frac{1}{2}(\underline{x}-\underline{\mu})^T \underline{\Sigma}^{-1} (\underline{x}-\underline{\mu})} = \underline{\mu}$

data statistics  $\frac{1}{N} \sum_{i=1}^N x_i$

Proof. ML minimizes  $-\log \prod_{i=1}^N P(x_i | \underline{\lambda}) = -\sum_{i=1}^N \log P(x_i | \underline{\lambda})$   
For exponential distributions this is  
 $F(\underline{\lambda}) = N \log Z(\underline{\lambda}) - \sum_{i=1}^N \underline{\lambda} \cdot \underline{\phi}(x_i)$

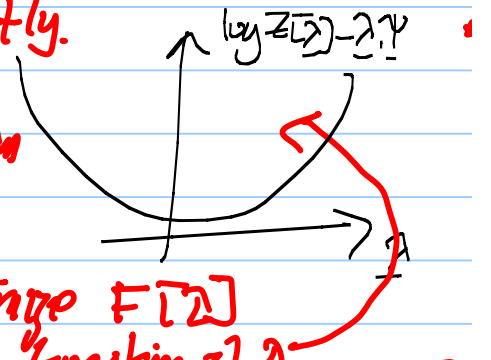
Differentiating w.r.t.  $\underline{\lambda}$

$$\frac{\partial F}{\partial \underline{\lambda}} = N \sum_{\underline{x}} \underline{\phi}(\underline{x}) P(\underline{x} | \underline{\lambda}) - \sum_{i=1}^N \underline{\phi}(x_i) \quad \text{result follows.}$$

Note: for some exponential distributions it is possible to compute  $\sum_{\underline{x}} \underline{\phi}(\underline{x}) P(\underline{x} | \underline{\lambda})$  analytically as a function of  $\underline{\lambda}$  (e.g. Gaussian), and hence solve ML directly.

But for other exponential distributions we cannot compute  $\sum_{\underline{x}} \underline{\phi}(\underline{x}) P(\underline{x} | \underline{\lambda})$ .

Instead we use an algorithm to minimize  $F(\underline{\lambda})$  w.r.t.  $\underline{\lambda}$ . Fortunately  $F(\underline{\lambda})$  is a convex function of  $\underline{\lambda}$  and hence has only a single minimum (proof  $\frac{\partial^2 F}{\partial \underline{\lambda}^2}$  is positive definite)



(7)

Algorithms for minimizing  $F(\lambda)$  include:

(i) Steepest Descent:

$$\lambda^{t+1} = \lambda^t - \Delta \frac{\partial F}{\partial \lambda}$$

$\Delta$  is a small constant.

$\rightarrow$  Newton-Raphson.

$$\lambda^{t+1} = \lambda^t - \Delta \left\{ \frac{\sum_{\underline{x}} \phi(\underline{x}) P(\underline{x}|\lambda^t) - \frac{1}{N} \sum_{i=1}^N \phi(\underline{x}_i)}{\sum_{\underline{x}} \phi(\underline{x}) P(\underline{x}|\lambda^t)} \right\}$$

(ii) Generalized Iterative Scaling (GIS)

$$\lambda^{t+1} = \lambda^t - \log \frac{\sum_{\underline{x}} \phi(\underline{x}) P(\underline{x}|\lambda^t)}{\frac{1}{N} \sum_{i=1}^N \phi(\underline{x}_i)}$$

Comment: Steepest descent requires specifying a step size  $\Delta$ . If  $\Delta$  is too big, algorithm fails to converge. If  $\Delta$  is too small, algorithm converges slowly. Both algorithms guaranteed to converge to the correct solution (provided  $\Delta$  is well-chosen).

Note: both algorithms require computing

$$\sum_{\underline{x}} \phi(\underline{x}) P(\underline{x}|\lambda^t) \text{ at each stage } t \text{ of the algorithm}$$

Computing this can be difficult to perform numerically for some distributions. If so, stochastic sampling methods like Markov Chain Monte Carlo (MCMC) may be used.

(2)

Examples of learning Exponential.

$$P(x|z) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Gaussian: Statistics  $\phi(x) = (x, x^2)$ , OR  $\phi(x) = (x, x^T)$

$$\sum_x P(x|z) (x, x^2) = \frac{1}{N} \sum_{i=1}^N (x_i, x_i^2) \quad \text{for Gaussian in } N\text{-Dim. } x \text{ } N\text{-D vector}$$

L.H.S.  $\int P(x|z) x = \mu$

$$\int P(x|z) x^2 = \mu^2 + \sigma^2,$$

Hence  $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$ ,  $\hat{\mu}^2 + \hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N x_i^2$ , so  $\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2$

Note: Really should be  $\int dx P(x|z)$  for Gaussian.

Letter Example:  $x \in A = \{a, b, c, d, e, \dots, y, z\}$

Statistic  $\phi(x) = (S_{x_a}, S_{x_b}, \dots, S_{x_z})$

ie.  $\phi(c) = (0, 0, 1, 0, 0, \dots, 0)$

$S_{x_a} = 1$ , if  $x = a$   
 $= 0$ , otherwise



(9)

Dataset of letters  $X_N = \{X_1, \dots, X_N\}$

$$\frac{1}{N} \sum_{i=1}^N \phi(x_i) = \left( \frac{\#a's}{N}, \frac{\#b's}{N}, \dots, \frac{\#z's}{N} \right)$$

where  $\#a's = \sum_{i=1}^N S_{X_i a} = \text{no. of } a's$   
in dataset.

$$\lambda = (\lambda_a, \dots, \lambda_z)$$

$$P(x|\lambda) = \frac{1}{Z(\lambda)} e^{\lambda_a S_{x a} + \dots + \lambda_z S_{x z}}$$

$$\sum_x P(x|\lambda) S_{x a} = \frac{1}{Z(\lambda)} e^{\lambda_a}$$

$$\text{wsto: } P(x=a|\lambda) = \frac{e^{\lambda_a}}{Z(\lambda)}$$

Hence solution

$$\frac{1}{Z(\lambda)} (e^{\lambda_a}, e^{\lambda_b}, \dots, e^{\lambda_z}) = \left( \frac{\#a's}{N}, \frac{\#b's}{N}, \dots, \frac{\#z's}{N} \right)$$

Hence  $\hat{\lambda}_a = \log \frac{\#a's}{N} - \log N$ ,  $\hat{\lambda}_b = \log \frac{\#b's}{N} - \log N$ ,  $\dots$ ,  $\hat{\lambda}_z = \log \frac{\#z's}{N} - \log N$ .

$$Z(\hat{\lambda}) = \frac{\#a's + \#b's + \dots + \#z's}{N} = 1$$

$$\frac{S_0}{Z(\lambda)} = \frac{e^{\lambda_a + \dots + \lambda_z}}{Z(\lambda)}$$

(1b) Topic III: Kullback-Leibler & Approximate Distributions

The theory of ML estimation in the statistics literature assumes that  $P(x|\underline{\lambda})$  is the correct model for the data.

We now discuss how to justify ML as an approximation if the data is generated by a different distribution.

First define the Kullback-Leibler divergence between two distributions  $P(x|\underline{\lambda})$  and  $f(x)$

$$D(f||p) = \sum_x f(x) \log \left\{ \frac{f(x)}{P(x|\underline{\lambda})} \right\}$$

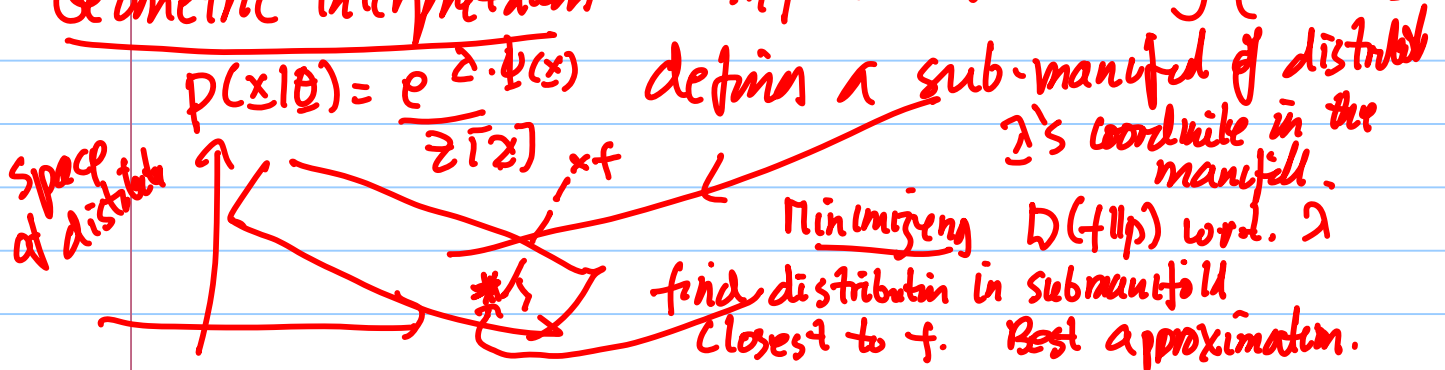
Properties:  $D(f||p) \geq 0$ ,  $D(f||p) = 0$  only if  $f(x) = P(x|\underline{\lambda})$   
i.e. if  $D(f||p)$  are small, then distributions  $f(x)$  and  $P(x|\underline{\lambda})$  are similar. If it is large, then they are not.

We can express:

$$D(f||p) = \underbrace{\sum_x f(x) \log f(x)}_{\text{index of } \underline{\lambda}} - \underbrace{\sum_x f(x) \log P(x|\underline{\lambda})}_{\text{function of } \underline{\lambda}}$$

Hence minimize  $D(f||p)$  w.r.t.  $\underline{\lambda}$  corresponds to minimizing  $-\sum_x f(x) \log P(x|\underline{\lambda})$

Geometric Interpretation: Information Geometry (Amari).



(11)

Relation to  $M_L$ .

Set  $f(x) = \frac{1}{N} \sum_{i=1}^N I(x=x_i)$  Indicator Function

This is the empirical distribution of the data  $\{x_i: i=1, \dots, N\}$   
(This is a special case of Parzen windows, later lecture.)

In this, minimizing the KL divergence corresponds to minimizing:

$$-\sum_x f(x) \log P(x|\lambda) = -\sum_x \frac{1}{N} \sum_{i=1}^N I(x=x_i) \log P(x|\lambda) = -\frac{1}{N} \sum_{i=1}^N \log P(x_i|\lambda).$$

This proves the

same criteria as ML!

Claim: ML estimation of  $\lambda$  is equivalent to minimizing  $D(f || P(x|\lambda))$  wrt.  $\lambda$ , where  $f(x)$  is the empirical distribution of the data.

Hence we can justify ML (for exponentials) as obtaining the distribution of form  $\frac{e^{\lambda \cdot \psi(x)}}{Z(\lambda)}$  which best approximates the data.

This also motivates the idea of model pursuit.

- (1) Start by doing ML on an exponential distribution with statistic  $\psi_1(x)$ . Get best approximation.
- (2) Get a better approximation by using more complex statistics  $\rightarrow$  e.g.  $\psi_1(x), \psi_2(x)$  with parameters  $\lambda_1, \lambda_2$ .
- (3) Proceed by using (increasingly) complex stats of course. Beyond scope

(12)

Example: Data consists of pairs of letters

$$X_N = (X_1^1, X_2^1), (X_1^2, X_2^2), \dots, (X_1^N, X_2^N)$$

First Model: assume independency

$$P(X_1, X_2) = P(X_1)P(X_2) \quad \text{with} \quad P(X) = \frac{1}{\sum T(x)} e^{-\lambda \cdot \phi(x)} \quad \text{as model for letters}$$

$$\text{Then} \quad \underline{P(X_1, X_2)} = \frac{e^{-\lambda \cdot (\phi(x_1) + \phi(x_2))}}{\sum T(x)} \quad (\lambda \text{ is estimated as before})$$

This gives best fit - in the Kullback-Leibler sense - to data, using statistic  $\phi_1(x_1, x_2) = \phi(x_1) + \phi(x_2)$ .

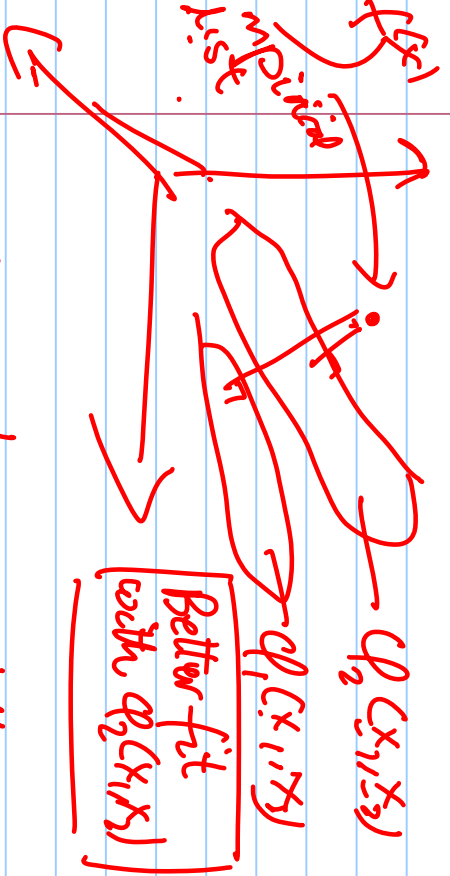
But we can use a better statistic:  $\phi_2(x_1, x_2)$

Pairwise frequencies of letters with  $\phi_2(x_1, x_2) = (S_{x_1 a} S_{x_2 a}, S_{x_1 b} S_{x_2 b}, S_{x_1 c} S_{x_2 c}, S_{x_1 d} S_{x_2 d}, S_{x_1 e} S_{x_2 e}, S_{x_1 f} S_{x_2 f}, S_{x_1 g} S_{x_2 g}, S_{x_1 h} S_{x_2 h}, S_{x_1 i} S_{x_2 i}, S_{x_1 j} S_{x_2 j}, S_{x_1 k} S_{x_2 k}, S_{x_1 l} S_{x_2 l}, S_{x_1 m} S_{x_2 m}, S_{x_1 n} S_{x_2 n}, S_{x_1 o} S_{x_2 o}, S_{x_1 p} S_{x_2 p}, S_{x_1 q} S_{x_2 q}, S_{x_1 r} S_{x_2 r}, S_{x_1 s} S_{x_2 s}, S_{x_1 t} S_{x_2 t}, S_{x_1 u} S_{x_2 u}, S_{x_1 v} S_{x_2 v}, S_{x_1 w} S_{x_2 w}, S_{x_1 x} S_{x_2 x}, S_{x_1 y} S_{x_2 y}, S_{x_1 z} S_{x_2 z})$

[13]

The second model, with  $q_2(x_1, x_2)$ , gives a better fit to the data - if the pairs of letters come from English.

Because there's pairwise regularization  $\rightarrow$  e.g.  $q_u$  is frequent  $q_z$  is impossible



Note: If we have more letters

e.g.  $X = \{ \text{brown, smith, loves, hates, ghost, ...} \}$  then higher order statistics are best.

But higher order statistics require more parameters -  $M$ -letters requires  $26^M$  parameters - so we don't usually have enough data. n.b. Shannon fit models like to data to estimate the entropy of English. See slide 16.

# (14) Maximum Entropy: Advanced Topic.

Another perspective on learning

Where do exponential distributions come from?  
Jaynes (wiki) claim they come from a maximum entropy principle.

We have data  $(x_1, \dots, x_n)$ .

We have statistics  $\phi(x)$  of the data.  
How to justify a distribution like  $P(x) = e^{-\lambda \cdot \phi(x)}$ ?

And how to justify using ML to get  $\lambda$ ?  $\frac{\partial \mathcal{H}(P)}{\partial \lambda}$

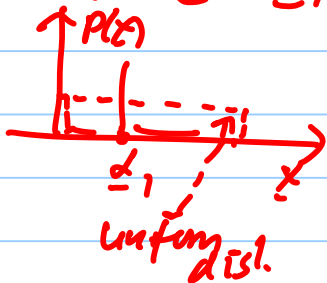
$$\mathcal{H}(P) = - \sum_x P(x) \log P(x) \quad \text{Entropy of Distribution.}$$

It is a measure of the amount of information that an observer expects to obtain by observing a sample  $x$  from a distribution  $P(x)$

$$\left. \begin{array}{l} \text{Info from sample} = -\log P(x) \\ \text{Expected info} = -\sum_x P(x) \log P(x) \end{array} \right\} \text{Shannon: Theory of Information}$$

Example: Suppose  $x$  can take  $M$  states  $\alpha_1, \dots, \alpha_M$

$$P(x = \alpha_1) = 1, \quad P(x = \alpha_i) = 0, \quad i = 2, \dots, M.$$



$$\text{Then } \mathcal{H}(P) = 0 \quad \left( \begin{array}{l} \text{Note } 0 \log 0 = 0 \\ 1 \log 1 = 0 \end{array} \right)$$

So no. info is gained by observing the sample, because we know it can only be  $\alpha_1$  (other possibilities can't happen).

Alternating. Suppose  $P(x = \alpha_i) = 1/M$ , for all  $i$ . Uniform Dist.  
 $\mathcal{H}(P) = \log M$ .

# (13) Maximum Entropy Principle

Given statistics  $\phi(x)$  with observed value  $\psi$ , choose the distribution  $P(x)$  to maximize the entropy subject to constraints

$$-\sum_x p(x) \log p(x) + \mu \left( \sum_x p(x) - 1 \right) + \lambda \cdot \left( \sum_x p(x) \phi(x) - \psi \right)$$

Lagrange multiplier      constraint      constraint

$\frac{\delta}{\delta p(x)}$

$$-\log p(x) - 1 + \mu + \lambda \cdot \phi(x) = 0$$

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

where  $\lambda, Z[\lambda]$  are chosen to satisfy the constraints:

$$\sum_x p(x) = 1, \Rightarrow Z[\lambda] = \sum_x e^{\lambda \cdot \phi(x)}$$

$$\sum_x p(x) \phi(x) = \psi, \Rightarrow \lambda \text{ is chosen s.t.}$$

$$\sum_x p(x|\lambda) \phi(x) = \psi$$

The maximum entropy principle recovers exponential distribution!

## (16) Entropy and ML.

Suppose we have data  $\{x_i: i=1 \text{ to } N\}$ , and fit a probability distribution  $P(x|\lambda)$  by ML - to get parameter  $\hat{\lambda}$

The prob of the data - using  $P(x|\hat{\lambda})$  <sup>← best estimate of  $\lambda$</sup>  is  $\prod_{i=1}^N P(x_i|\hat{\lambda}) = \exp\left\{\hat{\lambda} \sum_{i=1}^N \phi(x_i) - N \log Z(\hat{\lambda})\right\}$

The entropy of  $P(x|\hat{\lambda})$  is

$$\begin{aligned} -\sum_x P(x|\hat{\lambda}) \log P(x|\hat{\lambda}) &= \log Z(\hat{\lambda}) - \sum_x \hat{\lambda} \cdot \phi(x) P(x|\hat{\lambda}) \\ &= \log Z(\hat{\lambda}) - \frac{1}{N} \sum_{i=1}^N \hat{\lambda} \cdot \phi(x_i) \quad (\text{by actn. of } \hat{\lambda}) \end{aligned}$$

hence Prob of Data given  $P(x|\hat{\lambda})$   
 $= \exp\left\{-N \mathcal{L}[P(x|\hat{\lambda})]\right\} //$

So if the entropy of  $P(x|\hat{\lambda})$  is small, then it does not describe the data well - it cannot predict it and there is a lot of uncertainty.

Two Related means of the model. Its entropy.  
And the prob. of the data given the model.

Motivated Shannon search for the entropy of English.  
An example of model pursuit

Shannon starts with unary statistics - frequency of letters - fit to data (English text) estimated entropy.

then Shannon used more complex statistics - pairwise frequencies - entropy decreased (better fit)

and so on. Compare to human entropy  
what is the next letter to "ryth.?"



(17)

Which Models to Use?  
How to select between them?

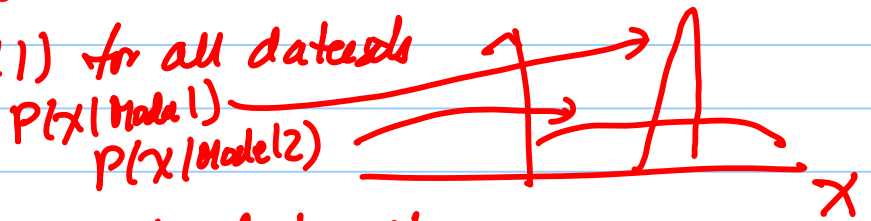
Suppose we have two models  $P(x | \text{Model 1})$ ,  $P(x | \text{Model 2})$   
e.g.  $P(x | \text{Model 1})$  uses stats  $\phi_1(x)$ , parameter  $\lambda_1$   
 $P(x | \text{Model 2})$  " "  $\phi_2(x)$ , parameter  $\lambda_2$ .

Let  $X = \{x_i : i \in \mathbb{N}\}$

Compute  $P(X | \text{Model 1})$  and  $P(X | \text{Model 2})$

Select Model 1 if  $P(X | \text{Model 1}) > P(X | \text{Model 2})$

Advantage - Occam Factor - this criterion favors simpler more specific models because the distribution must obey  $\sum_x P(X | \text{Model 1})$  for all datasets



Model 2 does okay on all datasets, but Model 1 does very well on some datasets and very badly on others. So if Model 1 does well on your dataset, then you should use it.

But there are complications. Model 1 corresponds to choice of statistics.  $\phi_1(x)$

$$P(X | \text{Model 1}) = \sum_{\lambda} P(X | \lambda) P(\lambda | \text{Model 1})$$

$$P(X | \lambda) = e^{-\lambda \cdot \phi(x)}$$

$$P(X | \lambda) = \prod_{i=1}^n P(x_i | \lambda)$$

Usually difficult, or impossible, to perform  $\sum_{\lambda}$ .

(18)

So try to approximate

$P(X | \text{Model 1})$  by  $\max_{\lambda} P(X | \lambda)$   
(assume  $P(\lambda | \text{Model 1})$  is uniform)

gives  $P(X | \text{Model 1}) \approx P(X | \hat{\lambda}) \rightarrow$  the ML prob,  
we discussed earlier, which relates to the entropy.

But now we no longer have  $\sum_{\lambda} P(\lambda | \text{Model 1}) = 1$  (even if we include the prior  $P(\hat{\lambda} | \text{Model 1})$ )

So Occam factor does not apply.

Need to penalize more complex models.

$\rightarrow$  e.g. Compare  $-\log P(X | \hat{\lambda}_1)$  model  $P(X | \lambda) = e^{\frac{\hat{\lambda}_1 \cdot \Phi(X)}{\Xi(\lambda_1)}}$   
with  $-\log P(X | \hat{\lambda}_2)$  model  $e^{\frac{\hat{\lambda}_2 \cdot \Phi_2(X)}{\Xi(\lambda_2)}}$

Suppose model 2 has more parameters

$\rightarrow$  e.g.  $\Phi_1(x) = (\phi_1^1(x), \dots, \phi_{10}^1(x))$ ,  $\lambda_1 = (\lambda_1^1, \dots, \lambda_{10}^1)$   
 $\Phi_2(x) = (\phi_1^2(x), \dots, \phi_{100}^2(x))$ ,  $\lambda_2 = (\lambda_1^2, \dots, \lambda_{100}^2)$

Model 2 has 100 parameters, Model 1 has 10

So Model 2 can adjust to data better.

Penalize more complex model  $\rightarrow$  AIC, BIC

Each model pays penalty.  $k$ ,  $n$ , or  $k_2 \log n$   
where  $n$  is the no. of parameters.

Other criteria  $\rightarrow$  Minimum Description Length (MDL)