

# Lecture 5

# Entropy and Model Selection

Note Title

1/26/2011

The entropy of a distribution  $P(x)$  is  
 $H[P] = -\sum_x P(x) \log P(x)$ .

It is a measure of the information obtained by observing a sample  $x$  from  $P(x)$ .

Shannon. Information Theory.  
 Physics. Thermodynamics / Stat Physics

Examples:

Let  $x \in \{\alpha_1, \dots, \alpha_N\}$ ,

ie.  $x$  has a finite set of values.

Case (i). Suppose  $P(x=\alpha_1)=1$ ,  $P(x=\alpha_i)=0$ ,  $i \neq 1$   
 There is no uncertainty, we know that the sample from  $P(x)$  will be  $\alpha_1$  before we observe it.

Entropy  $H[P] = -\sum_{i=1}^N P(x=\alpha_i) \log P(x=\alpha_i)$

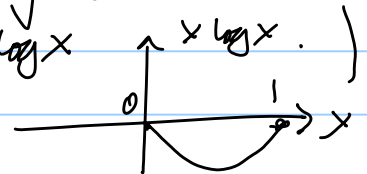
$$= -\{1 \log 1 + (N-1) 0 \log 0\}$$

$$= 0.$$

no information is gained by making the observation.

(Note:  $0 \log 0$  is defined by

$$\lim_{x \rightarrow 0} x \log x$$



Case (ii)  $P(x=\alpha_i) = 1/N$ ,  $i=1$  to  $N$ .

All samples are equally likely.

$$H[P] = -N \cdot (1/N) \log(1/N) = \log N.$$

we gain  $\log N$  bits of information by making the observation.

(2)

Note Title

1/26/2011

Shannon proposed that we encode information for transmission so that

events  $x$  will be encoded by  $-\log P(x)$  bits  
(i.e. high probability  $x$  have short codes,  
low probability  $x$  have long codes).

Then  $-\sum_x P(x) \log P(x)$  is the expected coding length of events  $x$  from distribution  $P(x)$ .

Shannon considered the task of modelling the English language

→ sequences of letters.

What is the entropy of English?

How does entropy relate to the learning ideas we discussed last time?

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

$$\begin{aligned} \text{Entropy} &= -\sum_x P(x|\lambda) \log P(x|\lambda) \\ &= \lambda \cdot \sum_x \phi(x) P(x) + \log Z(\lambda) \end{aligned}$$

Suppose we select  $\hat{\lambda}$  by minimizing the entropy such that  $\sum_x \phi(x) P(x) = \psi$  ← observed statistics

Then  $\hat{\lambda}$  is exactly the ML estimate  $\hat{\lambda}_{ML}$ .

(3)

### 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!

(4) An alternative viewpoint on ML learning of distributions. This gives deeper understand.

Suppose the data is generated by a distribution  $f(x)$ .

Define the Kullback-Leibler divergence between  $f(x)$  and the model  $P(x|\theta)$ .

~~Kullback-Leibler~~

$$D(f||P) = \sum_x f(x) \log \frac{f(x)}{P(x|\theta)}$$

KL has the property that

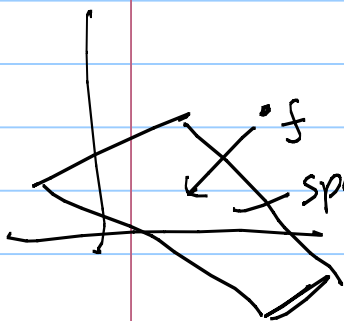
$$D(f||P) \geq 0 \quad \forall f, P.$$

$$D(f||P) = 0, \text{ if, and only if } f(x) = P(x|\theta).$$

So  $D(f||P)$  is a measure of the similarity between  $f(x)$  and  $P(x|\theta)$

We can write

$$D(f||P) = \underbrace{\sum_x f(x) \log f(x)}_{\text{Independent of } \theta} - \underbrace{\sum_x f(x) \log P(x|\theta)}_{\text{depends on } \theta}$$



space of all distributions  $p(x|\theta)$

Independent of  $\theta$ .

depends on  $\theta$

(5) Now suppose we have samples (i.i.d.)  $x_1, \dots, x_N$  from  $f(x)$ .

This gives us an empirical distribution  
empirical.  $\rightarrow f_{\text{emp}}(x) = \frac{1}{N} \sum_{i=1}^N \delta_{x, x_i}$   $\leftarrow$  Kronecker delta.  
Indicator Function.

The KL divergence between  $f_{\text{emp}}(x)$  and  $P(x|\theta)$  can be written as:

$$J(\theta) = - \sum_x f_{\text{emp}}(x) \log P(x|\theta) + K$$
$$J(\theta) = - \frac{1}{N} \sum_{i=1}^N \log P(x_i|\theta) + K$$

$\leftarrow$  independent of  $\theta$

Minimizing  $J(\theta)$  w.r.t.  $\theta$ , finds the distribution  $P(x|\hat{\theta})$  which is closest to  $f_{\text{emp}}(x)$ .

But minimizing  $J(\theta)$  w.r.t.  $\theta$  is exactly ML.  $\hat{\theta} = \text{ARG MAX}_{\theta} \sum_{i=1}^N \log P(x_i|\theta)$ . //

So ML has meaning even if best fit to the model. Even if model is only an approximation.

(6)

## Maximum Entropy + Model Selection

Two important papers (Della Pietra et al, Zhu et al) show how this principle is useful for learning complicated models.

(Della Pietra paper is available of the website.)

For a real problem, we do not know which statistics to use. (E.g., stats of letters, letter pairs etc.)

But we can define a large class of possible statistics.

$$\phi_1(x), \dots, \phi_c(x)$$

where  $c$  is very large

Two problems:

(i) selection: which of these statistics to use.

(ii) parameter estimation: how to determine the  $\lambda$  parameters of the model.

(7) Example: Della Pietra et al. (Zhu, Wu, Mumford)

Want a model to generate strings of text  $w = (w_1, w_2, \dots, w_n)$

Each  $w_i$  can be a letter, a space, or commas, etc.

first order model.

Simplest model: (lower case letters only)

$$p(w|N) = \frac{1}{Z} e^{\sum_{i=1}^N \lambda_{[a-z]} \cdot \chi_{[a-z]}(w_i)}$$

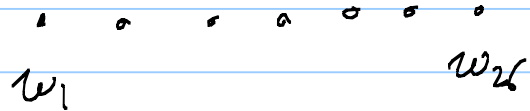
$\lambda$  parameters  
 $\chi$  potentials.  
 In this case indicator function.

length of string is generated by another process

$$\begin{aligned} \chi_{[a-z]}(w_i) &= (1, 0, \dots, 0), & \text{if } w_i = a \\ &= (0, 0, \dots, 0, 1), & \text{if } w_i = z \\ \lambda_{[a-z]} &= (\lambda_1, \lambda_2, \dots, \lambda_{26}), \end{aligned}$$

Represent as a graph.

Distribution is independent on letters



i.e.  $p(w|N) = \prod_{i=1}^N p(w_i)$  with  $p(w_i) = \frac{1}{Z} e^{\lambda_{[a-z]} \cdot \chi_{[a-z]}(w_i)}$

Learning the parameters corresponds to learning the frequencies of letters.

(8) This model does not fit the data (typical words). Because the letters are not independent - i.e. "t" is often followed by "h", "q" is almost followed by "u".

So second-order model using a new statistic.  $\chi_{[a-z][a-z]}(w_{ij})$  - indicator function of adjacent letters.

But there are many other features we could use. // Kronecker Delta Indicator Function.

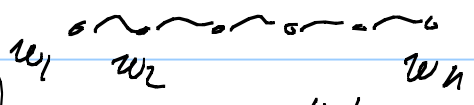
Stochastic sampling from the model is a generative way to check - model  
 samples from first order model xevu, ijjiir,  
 samples from second order model was, reaser, in,

Graphical Representation (dependent)

Suppose  $P(\underline{w}) = \frac{1}{Z(\underline{\lambda})} e^{\underline{\lambda} \cdot \underline{\phi}(\underline{w})}$  with  $\phi_{\mu}(\underline{w}) = \sum_{i=1}^{N-1} \delta_{f(w_{i+1}, w_i), \mu}$

$\sum_{\mu=1}^M \lambda_{\mu} \phi_{\mu}(\underline{w}) = \sum_{\mu=1}^M \sum_{i=1}^{N-1} \delta_{f(w_{i+1}, w_i), \mu} = \sum_{i=1}^{N-1} \lambda_{f(w_{i+1}, w_i)}$

$P(\underline{w}) = \frac{1}{Z(\underline{\lambda})} e^{\sum_{i=1}^{N-1} \lambda_{f(w_{i+1}, w_i)}}$  pairwise connection.





(9) The task of learning is to minimize  
 $\log Z[\underline{\lambda}] - \underline{\lambda} \cdot \underline{\psi}$  w.r.t.  $\underline{\lambda}$ .

Equivalently, minimize  $G[\underline{\lambda}] = Z[\underline{\lambda}] e^{-\underline{\lambda} \cdot \underline{\psi}}$

Write:  $G[\underline{\lambda}] = \sum_x e^{-\sum_{\mu} \lambda_{\mu} (\phi_{\mu}(x) - \psi_{\mu})}$  //

Initialize this by setting  $\lambda_{\mu} = 0, \forall \mu$ .  
(i.e. no statistics are selected)

Minimize  $G[\underline{\lambda}]$  by a "feature pursuit" strategy. (This will be used later in the course also).

At time  $t$ , we have state  $\{\lambda_{\mu}^t\}$  ( $\lambda_{\mu}^t = 0$ )

(A). Minimize w.r.t. each  $\lambda_v$  separately to solve for  $\lambda_v^{t+1}$  by solving  $\frac{\partial G[\underline{\lambda}]}{\partial \lambda_v} = 0$ , keep other  $\lambda_{\mu}$  fixed.

(B). Select the feature  $v$ , that minimizes  $G[\{\lambda_{\mu}^t: \mu \neq v, \lambda_v^{t+1}\}]$

(C) Updated  $\lambda_v^t \rightarrow \lambda_v^{t+1}$ .

Problem: both stages (A) & (B) are difficult

(10) New Page: Simply steps (A) & (B)

Minimize  $G(\underline{\lambda})$  w.r.t.  $\underline{\lambda}_v$ , with  $\underline{\lambda}_\mu$  fixed,  $\forall \mu \neq v$ .

$\underline{\lambda}_\mu = \underline{\lambda}_\mu^+$ ,  $\forall \mu \neq v$ ,  $\underline{\lambda}_v = \underline{\lambda}_v^+ + \underline{\Delta}_v$ , with  $\underline{\Delta}_v$  unknown.

$$\frac{\partial G}{\partial \underline{\lambda}_v} = \sum_x (\phi_v(x) - \psi_v) e^{-\sum_{\mu \neq v} \underline{\lambda}_\mu (\phi_\mu(x) - \psi_\mu)}$$

write as 
$$\frac{\partial G}{\partial \underline{\lambda}_v} = \sum_x (\phi_v(x) - \psi_v) e^{-\underline{\Delta}_v \cdot \phi_v(x)} e^{-\sum_{\mu \neq v} \underline{\lambda}_\mu^+ (\phi_\mu(x) - \psi_\mu)}$$

Hence  $\frac{\partial G}{\partial \underline{\lambda}_v} = 0$ , is equivalent to solving

We drop the constant  $e^{-\sum_{\mu \neq v} \underline{\lambda}_\mu^+ (\phi_\mu(x) - \psi_\mu)}$  and normalize.

(A) 
$$\sum_x (\phi_v(x) - \psi_v) e^{-\underline{\Delta}_v \cdot \phi_v(x)} p(x | \{\underline{\lambda}_\mu^+\}) = 0, \text{ for } \underline{\Delta}_v$$

We can use a similar technique to represent

$$\begin{aligned} G(\underline{\lambda}) &= \sum_x e^{-\sum_{\mu \neq v} \underline{\lambda}_\mu (\phi_\mu(x) - \psi_\mu)} \\ &= \sum_x e^{-\underline{\Delta}_v \cdot (\phi_v(x) - \psi_v)} e^{-\sum_{\mu \neq v} \underline{\lambda}_\mu^+ (\phi_\mu(x) - \psi_\mu)} \\ &= \sum_x e^{-\underline{\Delta}_v \cdot (\phi_v(x) - \psi_v)} e^{-\sum_{\mu \neq v} \underline{\lambda}_\mu^+ (\phi_\mu(x) - \psi_\mu)} \end{aligned}$$

where  $\underline{\lambda}_\mu = \underline{\lambda}_\mu^+$ , for all  $\mu \neq v$ ,  $\underline{\lambda}_v = \underline{\lambda}_v^+ + \underline{\Delta}_v$ .

Hence 
$$G(\underline{\lambda}) = \left( Z(\underline{\lambda}_\mu^+) e^{-\sum_{\mu \neq v} \underline{\lambda}_\mu^+ (\phi_\mu(x) - \psi_\mu)} \right) \sum_x e^{-\underline{\Delta}_v \cdot (\phi_v(x) - \psi_v)} p(x | \{\underline{\lambda}_\mu^+\})$$
  
 indep of  $v$ .

(B) To minimize  $G(\underline{\lambda})$  w.r.t.  $v$ , need to minimize  $\rightarrow \sum_x e^{-\underline{\Delta}_v \cdot (\phi_v(x) - \psi_v)} p(x | \{\underline{\lambda}_\mu^+\})$

(11) But there are some situations where this is not difficult. (see Della Pietra et al.)

Let the features  $\phi_\mu(\underline{x})$  be binary-valued.

Define:  $X_\nu^+ = \{ \underline{x} : \phi_\nu(\underline{x}) = +1 \}$

$X_\nu^- = \{ \underline{x} : \phi_\nu(\underline{x}) = -1 \}$ .

Express (1) as: (for this special case)

$$\sum_{\underline{x} \in X_\nu^+} (1 - \psi_\nu) e^{\Delta_\nu} p(\underline{x} | \langle \lambda_{-\mu}^+ \rangle) + \sum_{\underline{x} \in X_\nu^-} (-1 - \psi_\nu) e^{-\Delta_\nu} p(\underline{x} | \langle \lambda_{-\mu}^+ \rangle)$$

Hence  $\hat{\Delta}_\nu = \frac{1}{2} \log \left\{ \frac{(1 + \psi_\nu) \sum_{\underline{x} \in X_\nu^-} p(\underline{x} | \langle \lambda_{-\mu}^+ \rangle)}{(1 - \psi_\nu) \sum_{\underline{x} \in X_\nu^+} p(\underline{x} | \langle \lambda_{-\mu}^+ \rangle)} \right\}$

Express (2) as:

Select  $\nu$  to minimize:  $e^{\hat{\Delta}_\nu} \sum_{\underline{x} \in X_\nu^+} p(\underline{x} | \langle \lambda_{-\mu}^+ \rangle) + e^{-\hat{\Delta}_\nu} \sum_{\underline{x} \in X_\nu^-} p(\underline{x} | \langle \lambda_{-\mu}^+ \rangle)$

So we need to evaluate:  $\sum_{\underline{x} \in X_\nu^+} p(\underline{x} | \langle \lambda_{-\mu}^+ \rangle)$  &  $\sum_{\underline{x} \in X_\nu^-} p(\underline{x} | \langle \lambda_{-\mu}^+ \rangle)$

(12) It remains to compute:

(#)  $\sum_{\underline{x} \in \mathcal{X}_v^+} P(\underline{x} | \{\lambda_{\mu}^+\})$ ,  $\sum_{\underline{x} \in \mathcal{X}_v^-} P(\underline{x} | \{\lambda_{\mu}^-\})$   
(The sum of both of them is 1) for every  $v$ .

If the distribution  $P(\underline{x} | \{\lambda_{\mu}^+\})$  is defined over a finite set of states (Della Pietra et al) then these can be computed easily. //

Otherwise: Stochastic Sampling (eg. Jun Liu's Book, Stat 202C)

Key Point: Stochastic Sampling is very effective at estimating summations/integrals in high dimensions.

ie generate data from current model.

Also, for this problem, we only need to sample from  $P(\underline{x} | \{\lambda_{\mu}^+\})$  once to get sample  $\{\underline{x}_1, \dots, \underline{x}_M\}$   $M \gg 1$

Then we can compute (#) from these samples without needing to sample separately for each  $v$ . //

## (13) Penalizing Additional Features.

If our algorithm selects a feature  $v$  that is already non-zero, then we are merely adjusting the parameter  $\beta_v$  of that feature.

But if it selects a feature  $v$  with  $\beta_v = 0$ , then we are activating a feature that has not been used yet.

It can be good to penalize the no. of features used. There are a variety of criteria used to do this AIC, BIC,

$$kN \quad \text{or} \quad k \log N. \quad k - \text{constant.}$$

(14)

## Model Selection

To compare two models:

Model 1:

$$P(\underline{x} | \underline{\theta}, \text{Model 1}), P(\underline{\theta} | \text{Model 1})$$

$$P(\underline{x} | \underline{\phi}, \text{Model 2}), P(\underline{\phi} | \text{Model 2})$$

$$P(\underline{x} | \text{Model 1}) = \sum_{\underline{\theta}} P(\underline{x} | \underline{\theta}, \text{Model 1}) P(\underline{\theta} | \text{Model 1})$$

$$P(\underline{x} | \text{Model 2}) = \sum_{\underline{\phi}} P(\underline{x} | \underline{\phi}, \text{Model 2}) P(\underline{\phi} | \text{Model 2})$$

Then use the log-likelihood ratio test:

$$\log \frac{P(\underline{x} | \text{Model 1})}{P(\underline{x} | \text{Model 2})} \geq T \leftarrow \text{Threshold.}$$

Observe: that this takes into account all the different parameter values that could generate the data

Give an Occam factor — penalizes model with unnecessary large no. of parameters (\*)

(9)

## Minimax versus Model Selection

Zhu et al. proposed learning distributions by minimax.

(1) Use Maximum Entropy to learn a set of distributions using different statistics.

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

$\phi(x)$  are chosen statistics

$\psi = \sum_x \phi(x) P(x|\lambda)$   
observed statistics.

(2) Use Minimum Entropy to select between different distributions - i.e. to compare

$$P_1(x|\lambda_1) = \frac{1}{Z[\lambda_1]} e^{\lambda_1 \cdot \phi_1(x)} \quad \psi_1$$
$$P_2(x|\lambda_2) = \frac{1}{Z[\lambda_2]} e^{\lambda_2 \cdot \phi_2(x)} \quad \psi_2$$

Entropy is 
$$-\sum_x P(x|\lambda) \log P(x|\lambda)$$

$$= \log Z[\lambda] - \lambda \cdot \psi$$

(10) Hence using the minimum entropy principle is like maximum likelihood.

→ both involve minimizing  
 $\log Z[\underline{\lambda}] - \underline{\lambda} \cdot \underline{\psi}$

Consider:  $P(\underline{x} | \underline{\lambda}_1, \underline{\lambda}_2) = \frac{1}{Z[\underline{\lambda}_1, \underline{\lambda}_2]} e^{\underline{\lambda}_1 \cdot \frac{1}{2} \underline{\phi}(\underline{x}) + \underline{\lambda}_2 \cdot \underline{\phi}(\underline{x})}$

Using minimum entropy is choosing

$(\underline{\lambda}_1 = \underline{\lambda}_1, \underline{\lambda}_2 = 0)$  over  $(\underline{\lambda}_1 = 0, \underline{\lambda}_2 = 0)$

provided  $P(\underline{x} | \underline{\lambda}_1, 0) > P(\underline{x} | 0, \underline{\lambda}_2)$ .

This differs from model selection. Where you should compare

$$\log \frac{\sum_{\underline{\lambda}_1} P(\underline{x} | \underline{\lambda}_1) P(\underline{\lambda}_1)}{\sum_{\underline{\lambda}_2} P(\underline{x} | \underline{\lambda}_2) P(\underline{\lambda}_2)} //$$



(17) Occam Factor: Mackay.

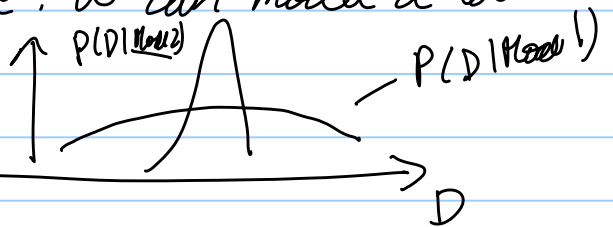
Automatic selection of complexity of the model.

$P(D | \text{Model 1})$

$P(D | \text{Model 2})$

Suppose Model 1 is flexible  $\rightarrow$  i.e. it can model a lot of datasets.

But Model 2 is more specific



Then Model selection will automatically favor the more specific model.

The Occam factor: