

Exponential Model  $P(d|\lambda) = \frac{e^{-\lambda} \lambda^{\phi(d)}}{Z[\lambda]}$

Data  $D = \{d^\mu; \mu=1 \dots N\}$   
 $P(D|\lambda) = \prod_{\mu} P(d^\mu|\lambda)$

Maximum Likelihood (ML)  $\hat{\lambda} = \text{ARG MAX}_{\lambda} P(D|\lambda)$   
 reduces to solve for  $\hat{\lambda}$   $\frac{1}{N} \sum_{\mu=1}^N \phi(d^\mu) = \sum_d \phi(d) P(d|\hat{\lambda})$

Probability of the data with best  $\lambda$

$$P(D|\hat{\lambda}) = \prod_{\mu} P(d^\mu|\hat{\lambda})$$

Intuitively, if  $P(D|\hat{\lambda})$  is big - i.e. the data is very probable  $\rightarrow$  then we think that the model fits the data well.

Model Selection:

Suppose we have two models

$$P_1(d|\lambda) = \frac{1}{Z_1[\lambda]} e^{\lambda \cdot \phi_1(d)}, \quad P_2(d|\lambda) = \frac{1}{Z_2[\lambda]} e^{\lambda \cdot \phi_2(d)}$$

with different statistics.

Which model is best?

For each model, find  $\hat{\lambda}_i = \text{ARG MAX}_{\lambda_i} P_i(D|\lambda_i)$   
 the best parameter for each model.

Evaluate:  $P_1(D|\hat{\lambda}_1) = \prod_{\mu=1}^N P_1(d_\mu|\hat{\lambda}_1)$

and  $P_2(D|\hat{\lambda}_2) = \prod_{\mu=1}^N P_2(d_\mu|\hat{\lambda}_2)$

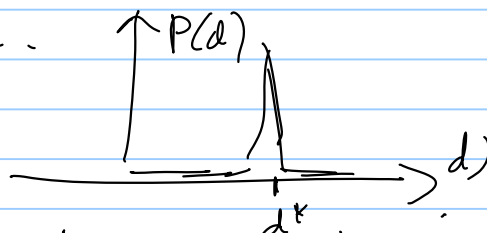
Select model 1, if  $P_1(D|\hat{\lambda}_1) > P_2(D|\hat{\lambda}_2)$   
 model 2, if  $P_2(D|\hat{\lambda}_2) > P_1(D|\hat{\lambda}_1)$   
 Model Selection (type I)

Note - there is an interpretation of this based on entropy.

Entropy  $H[P] = - \sum_d P(d) \log P(d)$

Entropy is a measure of how much information we gain from making an observation  $d$ .

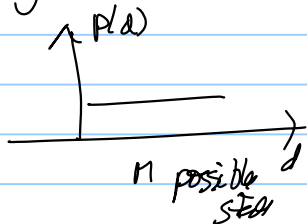
Suppose:  $P_0(d) = \delta(d-d^*)$   
 $= 0, \forall d \neq d^*$



We get no information from observing  $d^*$ , because we know that it is the only observation we can get.

$H(P_0) = 0$  ( $0 \log 0 = 0, 1 \log 1 = 0$ )

Suppose:  $P_1(d) = U(d)$  ← uniform distribution



then we get  $H(P_1) = \log M$  \* no. of possible values of  $d$

Result -  $\log P(D|\hat{\lambda}) = \sum_{\mu} \log P(d_{\mu}|\hat{\lambda})$   
 $\xrightarrow{\text{MLE estimate}} = \sum_{\mu} \hat{\lambda} \cdot \phi(d_{\mu}) - N \log Z[\hat{\lambda}]$

Entropy of  $P(d|\hat{\lambda})$   
 $= - \sum_d P(d|\hat{\lambda}) \log P(d|\hat{\lambda})$   
 $= - \sum_d P(d|\hat{\lambda}) \{ \hat{\lambda} \cdot \phi(d) - \log Z[\hat{\lambda}] \}$

But,  $\frac{1}{N} \sum_{\mu} \phi(d_{\mu}) = \sum_d \phi(d) P(d|\hat{\lambda})$   
 definition of ML

Hence  $\log P(D|\hat{\lambda}) = -N \text{Entropy}(P(d|\hat{\lambda}))$

$P(D|\hat{\lambda}) = e^{-N \text{Entropy}(P(d|\hat{\lambda}))}$

So  $P(D|\hat{\lambda})$  is big if Entropy  $P(d|\hat{\lambda})$  is small.

Maximum Likelihood corresponds to minimizing Entropy

The best model to fit data has lowest energy, hence best ability to predict.

Feature Pursuit:

Make a dictionary  $A = \{ \phi_1(d), \dots, \phi_n(d) \}$   
of possible features.

Task: want to construct a probability model.

$$P(d | \underline{\lambda}_{-i}) = \frac{1}{Z(\underline{\lambda}_{-i})} e^{\sum_{i=1}^n \lambda_i \phi_i(d)}$$

to model the data,

Want to keep model simple — use only a few of the features (also data limitation — later in course).

want  $\lambda_i = 0$  for most  $i$

Two tasks:

(i) selection — which features to use  
(i.e. to have  $\lambda_i \neq 0$ )

(ii) weighting — how to weight features and assign  $\lambda_i$ ?

This is a hard search problem (easier for discriminative learning — later in the course).

Strategy: feature Pursuit.  $\rightarrow$  Della Pietra<sup>12</sup>, Lafferty  
 $\rightarrow$  Zhu, Wu, Mumford.

(1) Find best model with one feature only

Calculate:  $\hat{i}$  s.t.  $P_n(D | \hat{\lambda}_{\hat{i}}) \geq P_i(D | \hat{\lambda}_i)$

Here  $P_i(d | \lambda_i) = \frac{1}{Z(\lambda_i)} e^{\lambda_i \phi_i(d)}$ ,  $\hat{\lambda}_i$  by ML for all  $i=1 \dots n$

This selects feature  $\phi_{\hat{i}}$  and assigns it weight  $\hat{\lambda}_{\hat{i}}$

(2) Next add another feature/statistic:

Consider all models of form:

$$P_{i,j}(d | \hat{\lambda}_i, \lambda_j) = \frac{1}{Z(\hat{\lambda}_i, \lambda_j)} e^{\hat{\lambda}_i \cdot \phi_{\hat{i}}(d) + \lambda_j \cdot \phi_j(d)}$$

$\uparrow$  feature selected already       $\uparrow$  new feat.

Page 4

Select the second feature  $\hat{j}$  by finding

$$P_{\hat{i}\hat{j}}(D | \hat{\lambda}_i, \hat{\lambda}_j) \geq P_{\hat{i}j}(D | \hat{\lambda}_i, \hat{\lambda}_j)$$

for all  $j = 1, \dots, n$

Proceed to select and weight the third, fourth, fifth, ... features and weight them.

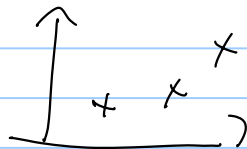
When to stop?

Adding a new feature allows the model to fit the data better

$$\rightarrow \text{i.e. } P_{\hat{i}\hat{j}}(D | \hat{\lambda}_i, \hat{\lambda}_j) \geq P_{\hat{i}}(D | \hat{\lambda}_i)$$

(because the model with two features has more flexibility and can find the data better)

- Example



fit data to line

$$y = a + bx$$

or to curve

$$y = a + bx + cx^2$$

easier,  $\rightarrow$  more flexibility

So, stop if increase by adding a new feature falls below a threshold:

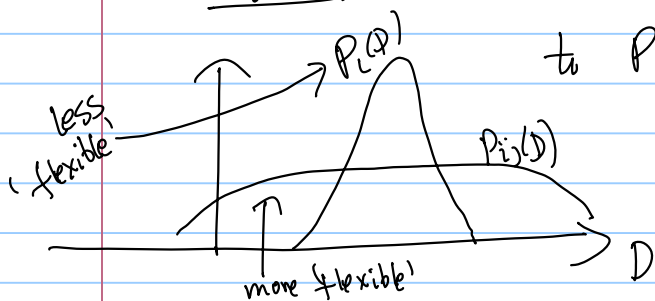
$\rightarrow$  i.e. stop with  $P_{\hat{i}}(D | \hat{\lambda}_i)$

$$\text{if } P_{\hat{i}\hat{j}}(D | \hat{\lambda}_i, \hat{\lambda}_j) \leq P_{\hat{i}}(D | \hat{\lambda}_i) + T \quad T = \text{threshold}$$

Note: A more advanced form of model selection will avoid this.  $\rightarrow$  Occam's razor (often impractical)

Require: Comparing  $P_i(D) = \sum_{\lambda_i} P_i(D | \lambda_i)$

$$\text{to } P_{ij}(D) = \sum_{\lambda_i, \lambda_j} P_{ij}(D | \lambda_i, \lambda_j)$$



$$\sum_D P_{ij}(D) = 1$$

$$\sum_D P_i(D) = 1$$

## Expectation-Maximization:

$$P(\underset{\text{observed}}{\underline{d}}, \underset{\text{hidden}}{\underline{h}} | \underline{\lambda}) = \frac{1}{Z(\underline{\lambda})} e^{\underline{\lambda} \cdot \phi(\underline{d}, \underline{h})}$$

Do ML on  $P(\underline{d} | \underline{\lambda}) = \sum_{\underline{h}} P(\underline{d}, \underline{h} | \underline{\lambda})$  to estimate  $\underline{\lambda}$

Minimize  $-\log P(\underline{d} | \underline{\lambda})$  with respect to  $\underline{\lambda}$

Add a new variable  $q(\underline{h})$ , a distribution over the hidden variables  $\sum_{\underline{h}} q(\underline{h}) = 1$

New  $P(\underline{h} | \underline{d}, \underline{\lambda})$  is the probability of the hidden variable  $\underline{h}$ , if we know the parameter  $\underline{\lambda}$ .

Formally  $P(\underline{h} | \underline{d}, \underline{\lambda}) = \frac{P(\underline{h}, \underline{d} | \underline{\lambda})}{\sum_{\underline{h}} P(\underline{h}, \underline{d} | \underline{\lambda})}$  (calculating  $\sum_{\underline{h}} P(\underline{h}, \underline{d} | \underline{\lambda})$  may be difficult - see later.)

Want  $q(\underline{h})$  to be close to  $P(\underline{h} | \underline{d}, \underline{\lambda})$

Kullback-Leibler  $\sum_{\underline{h}} q(\underline{h}) \log \frac{q(\underline{h})}{P(\underline{h} | \underline{d}, \underline{\lambda})} \geq 0$   
 $= 0$  if  $q(\underline{h}) = P(\underline{h} | \underline{d}, \underline{\lambda})$

Defn:

$$F(\underline{\lambda}, q(\cdot)) = -\log P(\underline{d} | \underline{\lambda}) + \sum_{\underline{h}} q(\underline{h}) \log \frac{q(\underline{h})}{P(\underline{h} | \underline{d}, \underline{\lambda})}$$

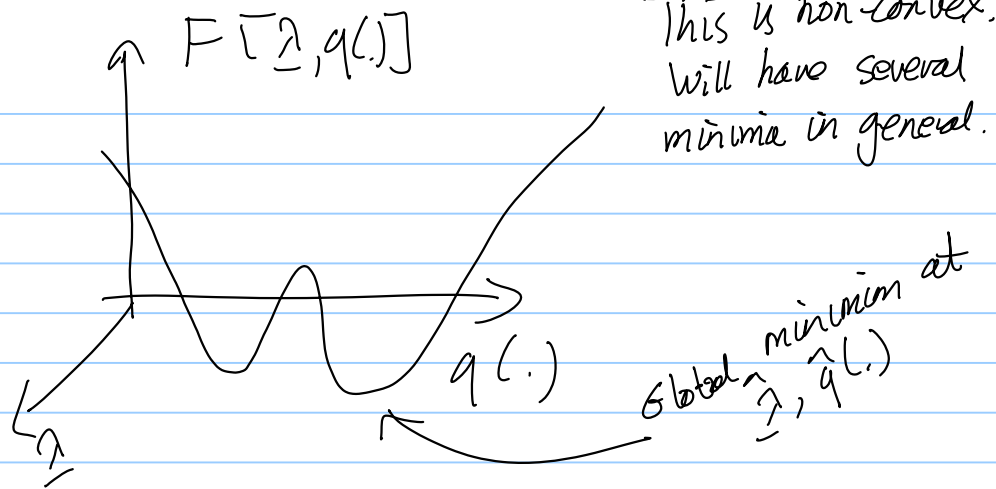
This is a function of  $\underline{\lambda}$  and  $q(\cdot)$ .

Its global minimum occurs at.

$$\hat{\underline{\lambda}} = \underset{\underline{\lambda}}{\text{ARG MIN}} (-\log P(\underline{d} | \underline{\lambda})), \quad \text{ML. } \hat{\underline{\lambda}}$$

and at  $\hat{q}(\underline{h}) = P(\underline{h} | \underline{d}, \hat{\underline{\lambda}})$  - makes second term 0.

Page 6.



Minimize by coordinate descent.

(1) At state  $\lambda^t$ ,

Solve for  $q^t(h) = \underset{q(\cdot)}{\text{ARG MIN}} F[\lambda^t, q(\cdot)]$

Solution:  $q^t(h) = P(h|d, \lambda^t)$   
 (requires computing  $\frac{P(h, d | \lambda^t)}{\sum_h P(h, d | \lambda^t)}$ )

(2) At state  $q^t(h)$

Compute  $\lambda^{t+1} = \underset{\lambda}{\text{ARG MIN}} F[\lambda, q^t]$

Solution:  $\sum_h q^t(h) \phi(\underline{d}, \underline{h}) = \sum_{h, d} \phi(\underline{d}, \underline{h}) P(h, d | \lambda)$

data statistics and expected  $\underline{h}$  w.r.t.  $q^t(h)$ 
model statistics

Repeat Steps.

Performance depends on the initial condition  $\lambda^{t_0}$ .

Previous notes (Tuesday) gave results for the extended case when we have data  $D = \{d_{\mu} : \mu = 1 \dots M\}$

Then replace  $F[\lambda, q(\cdot)]$  by

$$-\sum_{\mu} \log P(d_{\mu} | \lambda) + \sum_{\mu} \sum_{h_{\mu}} q_{\mu}(h_{\mu}) \log \frac{q_{\mu}(h_{\mu})}{P(h_{\mu} | d_{\mu}, \lambda)}$$