# Lecture 4

Note Title                                                                                    10/7/2006

Learning parametric probability models. This lecture deals with distributions like $P(\underline{x})$ or $P(\underline{x}, y)$. Later lectures will learn $P(y|\underline{x})$, the regression problem.

Parameterized model $P(X|\theta)$ ← model parameter

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

Assume that data is independent identically distributed (i.i.d).                    (product for independence).
$$P(X_1, .., X_N | \theta) = \prod_{i=1}^{N} P(X_i | \theta).$$

Choose: $\hat{\theta} = \text{ARG MAX } P(X_1..X_N|\theta) = \underset{\theta}{\text{ARG MAX }} \log P(X_1..X_N|\theta)$

Hence $P(X_1,..X_N(\hat{\theta}) \geqslant P(X_1 . X_N|\theta)$, for all $\theta$

$$\log P(X_1...X_N|\mu, \sigma) = \sum_{i=1}^{N} \log P(X_i|\mu, \sigma)$$

$$= -\sum_{i=1}^{N} \frac{(X_i - \mu)^2}{2\sigma^2} - \sum_{i=1}^{N} \log \sqrt{2\pi}\,\sigma.$$

Differentiate w.r.l. $\mu, \sigma$ gives
$$\frac{\partial}{\partial \mu} \log P(X_1..X_N|\mu, \sigma) = \frac{1}{\sigma^2} \sum_{i=1}^{N} (X_i - \mu).$$

$$\frac{\partial}{\partial \sigma} \log P(X_1..X_N|\mu, \sigma) = \frac{1}{\sigma^3} \sum_{i=1}^{N} (X_i - \mu)^2 - \frac{N}{\sigma}.$$

(2)

Maximia occurs at

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^{N} X_i$$

$$\hat{\sigma}_2 = \frac{1}{N} \sum_{i=1}^{N} (X_i - \hat{\mu})^2.$$

<u>Note</u>: The Gaussian is a special case. It is often impossible to solve $\frac{\partial}{\partial \theta} \log P(X_1 .. X_N | \theta) = 0$ analytically. An algorithm is required. (see later)

Now $\hat{\theta} = $ ARG MAX $P(X_1 .. X_N | \theta)$ is the ML estimator.

We could use the MAP estimator

$$\hat{\theta} = \text{ARG MAX } P(\theta | X_1, ... X_N)$$

$$P(\theta | X_1 .. X_N) = \frac{P(X_1 \quad X_N | \theta) P(\theta)}{P(X_1 .. X_N)}$$

$$\log P(\theta | X_1 .. X_N) = \underbrace{\sum_{i=1}^{N} \log P(X_i | \theta)}_{N \text{ terms}} + \underbrace{\log P(\theta)}_{1 \text{ term}} - \log P(X_1 .. X_N)$$

Note if N is large then the prior $P(\theta)$ usually has little effect.

We could use other estimators. We could formulate this in terms of Bayes risk and use a loss function which penalizes errors in $\theta$.

But this lecture will stick with ML estimation.

## Exponential Distributions

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

normalization factor.

$\underline{\lambda}$ – parameters

$\underline{\phi}(x)$ – statistics.

$$\underline{\lambda} = (\lambda_1, \lambda_2, \ldots, \lambda_M)$$

$$\underline{\phi}(x) = (\phi_1(x), \phi_2(x), \ldots, \phi_M(x))$$

<u>Almost</u> every named distribution can be <u>expressed</u> as an exponential distribution (particularly if we allow missing/hidden variables).

For Gaussian in 1-dimension.

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

$$P(x|\lambda) = \frac{1}{Z[\underline{\lambda}]} e^{\lambda_1 x + \lambda_2 x^2} \quad \text{compare to} \quad \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}.$$

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

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

(4)

## Learning an Exponential Distribution

You can learn them by Maximum Likelihood,

Examples:

$$P(\langle \underline{x}_1, \underline{x}_2, .., \underline{x}_N \rangle | \underline{\lambda}) = \prod_{\mu=1}^{N} e^{\frac{-\underline{\lambda} \cdot \underline{\phi}(\underline{x}_\mu)}{Z[\underline{\lambda}]}}$$

Maximize wrt $\underline{\lambda}$ !!

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

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

Claim: $\dfrac{\partial \log Z[\underline{\lambda}]}{\partial \underline{\lambda}} = \sum_{x} \underline{\phi}(x) P(\underline{x}|\underline{\lambda})$

expected value of the statistics $\underline{\phi}(x)$ w.r.t. $P(\underline{x}|\underline{\lambda})$

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

(5)

Claim: For exponential distribution, ML
corresponds to finding the value of $\lambda$ s.t. the
model statistics are equal to the data statistics

solve $\quad \sum_{\underline{x}} \phi(\underline{x}) P(\underline{x}|\underline{\lambda}) = \frac{1}{N} \sum_{i=1}^{N} \phi(\underline{x}_i)$ .

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 $\frac{1}{N} \sum_{i=1}^{N} \underline{x}_i$ .

Proof. ML minimizes $-\log \prod_{i=1}^{N} P(\underline{x}_i|\underline{\lambda}) = -\sum_{i=1}^{N} \log P(\underline{x}_i|\underline{\lambda})$

For exponential distributions this is,

$$F[\underline{\lambda}] = N \log Z[\underline{\lambda}] - \sum_{i=1}^{N} \underline{\lambda} \cdot \phi(\underline{x}_i)$$

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

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

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

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

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) $\frac{\partial^2 F}{\partial\underline{\lambda}\partial\underline{\lambda}}$ is positive definite)

$\log Z[\underline{\lambda}] - \underline{\lambda} \cdot \underline{\mu}$

(6)

Algorithms for minimizing $F[\underline{\lambda}]$ include:

(i) Steepest Descent:
$$\underline{\lambda}^{t+1} = \underline{\lambda}^t - \Delta \frac{\partial F}{\partial \underline{\lambda}}$$

$\Delta$ is a small constant

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

$\rightarrow$ Newton Raphson.

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

$\left( \text{Note} \quad \log \underset{\nearrow}{\underline{\psi}} = (\log \psi_1, \log \psi_2 \ldots, \log \psi_m) \quad, \text{ where } \underline{\psi} = (\psi_1 \ldots \psi_m) \right.$

vector | all components are positive.]

**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).

<u>Note</u>: both algorithms require computing
$$\sum_{\underline{x}} \underline{\phi}(\underline{x}) P(\underline{x}|\underline{\lambda}^t) \quad \text{at each stage } t \text{ of the algorithm}$$

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.

(7)

The theory of ML estimation in the statistics literature assumes that $P(x | \lambda)$ 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 $P(\underline{x} | \lambda)$ and $f(\underline{x})$

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

Properties: $D(f \| p) > 0$, $D(f \| p) = 0$ only if $f(\underline{x}) = P(\underline{x} | \lambda)$

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

We can express:

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

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

Geometric Interpretation : Information Geometry (Amari).

$$P(x | \theta) = \frac{e^{\lambda \cdot \underline{\phi}(x)}}{\overline{z}[\lambda]}$$

defines a sub-manifold of distributions

$\lambda$'s coordinate in the manifold.

Minimizing $D(f \| p)$ w.r.t. $\lambda$ find distribution in submanifold closest to $f$. Best approximation.

space of distrib.

$\times f$

$\ast \lambda$

(8)

## Relation to ML.

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

This is the empirical distribution of the data $\{\underline{x_i} : i = 1...N\}$

(This is a special case of Parzen windows, later lecture.)

In this, minimizing the KL divergence corresponds to minimizing:

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

This proves the

→ Same criteria as ML!

**Claim:** ML estimation of $\underline{\lambda}$ is equivalent to minimize $D(f || P(\underline{x}|\underline{\lambda}))$ w.r.t. $\underline{\lambda}$, where $f(\underline{x})$ is the empirical distribution of the data.

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

This also motivates the idea of _model pursuit_.

(1) Start by doing ML on an exponential distribution with statistic $\underline{\phi_1}(\underline{x})$. Get best approximation

(2) Get a better approximation by using more complex statistics → e.g. $\underline{\phi_1}(\underline{x}), \underline{\phi_2}(\underline{x})$ with parameters $\underline{\lambda_1}, \underline{\lambda_2}$

(3) Proceed by using increasingly complex stats of course.
Beyond scope

(9)

## Another perspective on learning

Where do exponential distributions come from?
Jaynes (wiki) claim they come from a __maximum__
__entropy principle__.

We have data $(x_1, \ldots x_N)$.
We have statistics $\phi(x)$ of the data,
How to justify a distribution like $P(x) = \dfrac{e^{\eta \cdot \phi(x)}}{z(\eta)}$?
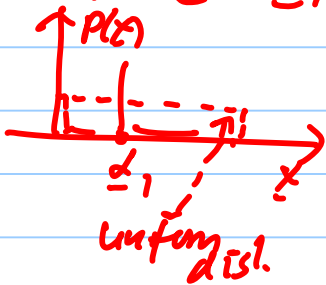And how to justify using ML to get $\eta$?

$$\mathcal{H}[p] = -\sum_x P(x) \log P(x) \qquad \text{Entropy of Distribution}$$

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

$$\text{Info from sample} \sim -\log P(x)$$
$$\text{Expected info} = -\sum_x P(x) \log P(x)$$

$\left.\begin{array}{c} \\ \\ \end{array}\right\}$ Shannon: Theory Information

__Example__: Suppose $x$ can take M states $\alpha_1, \ldots \alpha_M$

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

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

[sketch of $P(x)$ vs $x$ with spike at $\alpha_1$, labeled "uniform dist."]

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

Alternative. Suppose $P(x = \alpha_i) = 1/M$, for all $i$. Uniform Distr.
$\mathcal{H}[P] = \log M$.

**Maximum Entropy Principle**

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

Lagrange multipliers ↘   ← constraint

$$-\sum_{\underline{x}} P(\underline{x}) \log P(\underline{x}) + \mu \left\{ \sum_{\underline{x}} P(\underline{x}) - 1 \right\}$$

constraint ↗

Lagrange multiplier ⟶ $+ \underline{\lambda} \cdot \left\{ \sum_{\underline{x}} P(\underline{x}) \phi(\underline{x}) - \underline{\Psi} \right\}$

$$\frac{\delta}{\delta P(\underline{x})} \qquad -\log P(\underline{x}) - 1 + \mu + \underline{\lambda} \cdot \underline{\phi}(\underline{x}) = 0$$

Solution. $\quad P(\underline{x} \mid \underline{\lambda}) = \dfrac{e^{\underline{\lambda} \cdot \underline{\phi}(\underline{x})}}{Z[\underline{\lambda}]}$

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

$$\sum_{\underline{x}} P(\underline{x}) = 1 \quad , \quad \Rightarrow \quad Z[\underline{\lambda}] = \sum_{\underline{x}} e^{\underline{\lambda} \cdot \underline{\phi}(\underline{x})}$$

$$\sum_{\underline{x}} P(\underline{x}) \phi(\underline{x}) = \underline{\Psi} \quad , \quad \Rightarrow \quad \underline{\lambda} \text{ is chosen s.t.}$$

$$\sum_{\underline{x}} P(\underline{x} \mid \underline{\lambda}) \phi(\underline{x}) = \underline{\lambda}$$

The maximum entropy principle recovers exponential distribution!

(11)    **Entropy and M.L.**    Spring 2023.

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

The prob of the data = using $P(x \mid \hat{\lambda})$  ← best estimate of $\hat{\lambda}$.

is $\prod_{i=1}^{N} P(x_i \mid \hat{\lambda}) = \exp\left\{ \hat{\lambda} \cdot \sum_{i=1}^{N} \phi(x_i) - N \log Z[\hat{\lambda}] \right\}$

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

$-\sum_{x} P(x \mid \hat{\lambda}) \log P(x \mid \hat{\lambda}) = \log Z[\hat{\lambda}] - \sum_{x} \hat{\lambda} \cdot \phi(x) P(x \mid \hat{\lambda})$

$\qquad\qquad = \log Z[\hat{\lambda}] - \frac{1}{N} \sum_{i=1}^{N} \hat{\lambda} \cdot \phi(x_i)$

$\qquad\qquad\qquad\qquad$ (By detn. of $\hat{\lambda}$)

Hence    Prob of Data given $P(x \mid \hat{\lambda})$

$\qquad = \exp\left\{ -N \, \mathcal{H}[P(x \mid \hat{\lambda})] \right\}$ . //

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

$\underline{\text{Two Related measures}}$ 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

$\underline{\text{Shannon}}$ 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."?

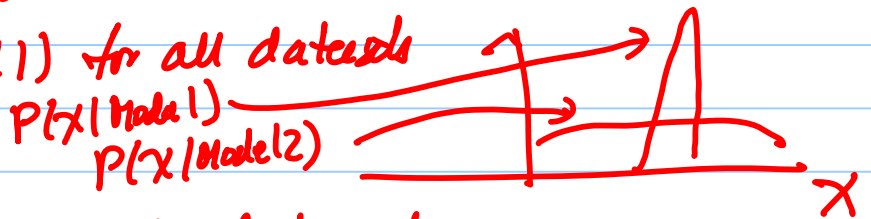(12)

## Which Models to Use?
### How to select between them?

Suppose we have two models $P(\underline{x}|Model 1)$, $P(\underline{x}|Model 2)$

e.g. $P(\underline{x}|Model 1)$ uses stats $\underline{\phi}_1(\underline{x})$, parameter $\underline{\lambda}_1$
$P(\underline{x}|Model 2)$  "    "  $\underline{\phi}_2(\underline{x})$, parameter $\underline{\lambda}_2$.

Let $X = \langle \underline{x}_i : i \in W \rangle$

Compute $P(X|Model 1)$ and $P(X|Model 2)$

Select Model 1 if $P(X|Model 1) > P(X|Model 2)$

Advantage — Occam Factor — this criterion favours simpler more specific models because the distribution must obey $\sum_X P(X|Model 1)$ for all datasets

$P(X|Model 1)$
$P(X|Model 2)$



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. $\phi_1(\underline{x})$

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

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

$$P(X|\underline{\lambda}) = \prod_{i=1}^{N} P(\underline{x}_i|\underline{\lambda}).$$

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

(13)

So try to approximate

$P(X \mid \text{Model 1})$ by $\max\limits_{\underline{\lambda}} P(X \mid \underline{\lambda})$

(assume $P(\underline{\lambda} \mid \text{Model 1})$ is uniform)

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

But now we no longer have $\sum\limits_X P(X \mid \text{Mod 1}) = 1$ (even if we include the prior $P(\hat{\underline{\lambda}} \mid \text{Model}))$

So Occam factor does not apply. //

Need to penalize more complex models.

$\to$ e.g. Compare $-\log P(X \mid \hat{\underline{\lambda}}_1)$ model $P(X \mid \underline{\lambda}) = e^{\dfrac{\hat{\underline{\lambda}}_1 \cdot \phi_1(x)}{Z[\underline{\lambda}_1]}}$

with $-\log P(X \mid \hat{\underline{\lambda}}_2)$ model $\dfrac{e^{\hat{\underline{\lambda}}_2 \cdot \phi_2(x)}}{Z[\underline{\lambda}_2]}$

Suppose model 2 has more parameters

$\to$ e.g. $\phi_1(\underline{x}) = (\phi_1^1(\underline{x}) \dots \phi_{10}^1(\underline{x})), \ \underline{\lambda}_1 = (\lambda_1^1, \dots \lambda_{10}^1)$

$\phi_2(\underline{x}) = (\phi_1^2(\underline{x}), \dots, \phi_{100}^2(\underline{x})), \ \underline{\lambda}_2 = (\lambda_1^2, \dots, \lambda_{100}^2)$

Model 2 has 100 parameters, Model 1 has 10

So Model 2 can adjust to data better.

Penalize more Complex model $\to$ AIC, BIC

Each model pays penalty. $k_1 n$, or $k_2 \log n$

where $n$ is the no. of parameters.

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

# Training and Testing

Critical issue is — how much data do you need to learn a distribution?

There is no perfect answer. A rule of thumb is that you need $k \times$ no. of parameters of the distribution, where $k = 5$ to $10$.

In practice, train (learn) the model on a training dataset. Test it on a second dataset. If performance (e.g. Bayes Risk, ROC curves, etc.) is the same on both — then you have learned or generalized.

If performance is good on the training set but bad on the testing set — then you have only memorized the training set.