

Exponential Distributions, Sufficient Statistics,
and the Maximum Entropy Principle.

Most distributions can be expressed in exponential form:

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

$\phi(\cdot)$ is the 'sufficient statistics'.

This can be derived from the maximum entropy principle:

ψ - derived from observation
see below

maximize $-\sum_d P(d) \log P(d) + \lambda \cdot \left(\sum_d P(d) \phi(d) - \psi \right)$

"entropy" Lagrange multiplier = expected stat - observed stat ψ

Differentiating w.r.t. $P(d)$ gives result.

The value of λ is chosen s.t. $\sum_d P(d) \phi(d) = \psi$.

Maximum Likelihood Estimation (MLE)

Data: $\langle d^N : \mu = H_0 \text{ or } N \rangle$

$$P(\langle d^N \rangle | \lambda) = \prod_{\mu} P(d^{\mu} | \lambda)$$

assumes data is independently identically distributed (i.i.d)

MLE

$$\hat{\lambda} = \text{Arg Max}_{\lambda} \left\{ -\log P(\langle d^N \rangle | \lambda) \right\}$$

$$= \text{Arg Max}_{\lambda} \left\{ N \log Z(\lambda) - \sum_{\mu=1}^N \lambda \cdot \phi(d^{\mu}) \right\}$$

$$= \text{Arg Max}_{\lambda} \left\{ \log Z(\lambda) - \lambda \cdot \psi \right\} \quad \text{where } \psi = \frac{1}{N} \sum_{\mu=1}^N \phi(d^{\mu})$$

Let. $G(\lambda) = \log Z(\lambda) - \lambda \cdot \psi$.

remark. $\frac{\partial}{\partial \lambda} \log Z(\lambda) = \sum_d \frac{\phi(d) e^{\lambda \cdot \phi(d)}}{Z(\lambda) e^{\lambda \cdot \phi(d)}} = \sum_d \phi(d) P(d|\lambda)$

$\frac{\partial^2}{\partial \lambda^2} \log Z(\lambda)$ is positive definite, hence $\log Z(\lambda) - \lambda \cdot \psi$ is convex, unique minimum.

minimum occurs when $\frac{\partial G}{\partial \lambda} = 0$

implies

$$\sum_d \phi(d) P(d|\lambda) = \psi$$

like maximum entropy

(2) Note: since $G[\lambda]$ is convex it may seem
easy to find its global minimum.

But (almost) all algorithms to find it will
require computing the gradient

→ eg. steepest descent

$$\lambda^{t+1} = \lambda^t - \epsilon \left(\sum_d \phi(d) P(d|\lambda^t) + \Psi \right)$$

this requires computing the expectation $\sum_d \phi(d) P(d|\lambda^t)$
which is difficult because computing the partition function
is very hard — except on trees.

Alternative algorithm

Generalized Iterative Scaling (GIS)

is a Discrete Iterative Algorithm (DIA).

~ simple modification of steepest descent, but guaranteed
to converge (without needing a step size ϵ).

$$\lambda^{t+1} = \lambda^t - \log \left(\sum_d \phi(d) P(d|\lambda^t) \right) + \log \Psi$$

(Note: \log is taken by components, e.g. $\log(a, b) = (\log a, \log b)$)

Note: we can always use MCMC to evaluate
 $\sum_d \phi(d) P(d|\lambda^t)$ → doesn't need to know $Z[\lambda]$.

Note: suppose we initialize $\lambda^0 = 0$
then $P(d|\lambda^0) = U(d)$ uniform distribution.

$$\lambda^1 = \log \Psi - \log \sum_d U(d) \phi(d)$$

For some applications: (i) $\sum_d U(d) \phi(d)$ can be calculated
(ii) λ^1 is a good estimate
i.e. one iteration of GIS gives
a good approximation

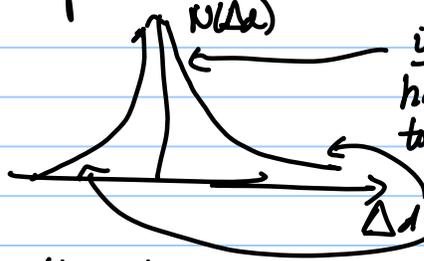
What $\phi(\cdot)$'s to use?

Claim: we can obtain MRF models from stats $\phi(d)$
in particular,
we can derive the "weak smoothness"

model on MRF's from this approach.

(3)

Claim: on most images the histogram of the difference operator $d_{i+1}-d_i$ takes the standard form



ie. most pixels have intensity similar to their neighbors, but there are some 'outliers', edges..

Note: this histogram is not Gaussian

A Gaussian would be too flat at $\Delta d = 0$ and fall away too quickly for large Δd ~ doesn't allow for outliers (Gaussians are not "robust").

statistic $\phi(z:d) = \frac{1}{N} \sum_{i=1}^N \delta(d_{i+1}-d_i, z)$ histogram

$$\lambda \cdot \phi(d) = \sum_z \lambda(z) \frac{1}{N} \sum_{i=1}^N \delta(d_{i+1}-d_i, z)$$

$$= \frac{1}{N} \sum_{i=1}^N \left\{ \sum_z \lambda(z) \delta(d_{i+1}-d_i, z) \right\}$$

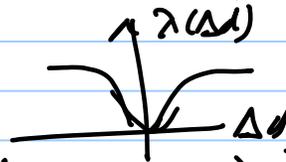
$$= \frac{1}{N} \sum_{i=1}^N \lambda(d_{i+1}-d_i) = \frac{1}{N} \sum_{i=1}^N \lambda(d_{i+1}-d_i)$$

distribution

$$p(d|\lambda) = \frac{1}{Z(\lambda)} e^{-\frac{1}{N} \sum_{i=1}^N \lambda(d_{i+1}-d_i)}$$

hence an MRF, with potentials $\lambda(\cdot)$ relating neighboring pixels.

How does "weak smoothness" relate to the histogram above?



Answer → one iteration of GIS (previous page) is enough.

What other statistics can we use?

Claim: any 'derivative operator' $\sum_{i=1}^m A_i d_i$ st. $\sum_{i=1}^m A_i$ gives similar statistics (Green)

Applying max entropy / exponential distribution will give MRFs with longer range interactions and higher order cliques.

Statistics can be combined

$$p(d|\lambda_1, \lambda_2) = \frac{1}{Z(\lambda_1, \lambda_2)} e^{\lambda_1 \phi_1(d) + \lambda_2 \phi_2(d)}$$

How to select between different combinations?

Model selection → prefer $p(d|\lambda_1, \lambda_2)$ to $p(d|\lambda)$ if $p(d|\lambda_1, \lambda_2) > p(d|\lambda)$ for some prior $p(\lambda_1, \lambda_2) > p(\lambda)$

Greedy search

Zhu, Wu, Mumford, Delta Pictorial

(4) Compare MLE to new machine learning alternatives.

First address simple regression problem

$$P(\omega | d) = \frac{e^{\omega \sum_{\mu} \lambda_{\mu} \phi_{\mu}(d)}}}{Z(d, \lambda)}$$

$$\omega \in \{+1\}, \text{ implies } Z(d, \lambda) = e^{\sum_{\mu} \lambda_{\mu} \phi_{\mu}(d)} + e^{-\sum_{\mu} \lambda_{\mu} \phi_{\mu}(d)}$$

Conditional distribution

Statisticians → treat this as a regression problem.

data $\{(w_i, d_i)\}$, maximize $\prod_i P(w_i | d_i)$

Note: regression was developed by Gauss around 1800 to predict the motion of a planetoid Ceres.

(regression can be applied to continuous & discrete variables).

leads to MLE problem:

$$\underset{\omega, \lambda}{\text{minimize}} \quad \sum_i \log Z(d_i; \lambda) - \sum_i \omega_i \lambda \cdot \phi(d_i)$$

$$\sum_i \log \left(e^{\lambda \cdot \phi(d_i)} + e^{-\lambda \cdot \phi(d_i)} \right) - \lambda \cdot \sum_i \omega_i \phi(d_i)$$

Note → this computation is practical since $Z(d_i; \lambda)$ can be evaluated.

Note: (i) statisticians usually work with regression models with a limited number of statistics $\phi_{\mu}(d)$.

(ii) statisticians usually minimize the function by steepest descent methods → update all the parameters λ simultaneously.

Contrast with machine learning (AdaBoost)

Changes: (i) call the statistics 'weak classifiers' and require them to be binary-valued.

(ii) minimize a related cost function

$$C_{\text{Ada}}(\lambda) = \sum_i e^{\omega_i \sum_{\mu} \lambda_{\mu} \phi_{\mu}(d_i)} \rightarrow \text{bound the classification loss}$$

(iii) algorithm: initialize $\lambda_{\mu} = 0$, fix

(other fixed) → coordinate descent. At each time step, pick the λ_{μ} which maximizes the decrease in cost.

(5) AdaBoost Algorithm

Initialize $\lambda_\mu = 0, \forall \mu.$

At time step t : $\{\lambda_\mu^t\}$

- For each μ , solve $\frac{\partial C_{\text{Ada}}(\lambda)}{\partial \lambda_\mu} = 0$

to solve for $\hat{\lambda}_\mu$

• Calculate $\hat{\mu} = \underset{\mu}{\text{ARG MIN}} C_{\text{Ada}}(\hat{\lambda} + \hat{\lambda}_\mu)$

• Set $\lambda_{\hat{\mu}}^{t+1} = \hat{\lambda}_{\hat{\mu}}^t$ $(\hat{\lambda} + \hat{\lambda}_\mu = (\lambda_1^t, \dots, \hat{\lambda}_\mu^t, \dots, \lambda_N^t))$
 $\lambda_\nu^{t+1} = \lambda_\nu^t$ for $\nu \neq \hat{\mu}$

Note: At time t , the algorithm either "selects" a new weak classifier $\phi_\mu(\cdot)$ (sh. $\lambda_\mu^t = 0$) and assigns it a non-zero "weight" λ_μ^t , or it changes the weight of a weak classifier that has already been selected.

→ AdaBoost (ie, Multiclass learning) puts much more emphasis on 'selection' → ie. you start with a dictionary of many weak classifiers $\{\phi_\mu(\cdot) : \mu \in \mathcal{M}\}$ and only select and use a small number (but small may mean 100's or 1,000's).

→ AdaBoost (Multiclass learning) emphasizes 'discrimination' rather than learning the conditional distribution $P(w|\mathbf{x})$
→ pays attention to data near the decision boundary and the margins.

→ arguably better than MLE if you have limited data
→ need to cross-validate to prevent over-learning.

→ The two algorithmic steps of AdaBoost

- ie. solve $\frac{\partial C_{\text{Ada}}}{\partial \lambda_\mu} = 0$, and $\hat{\mu} = \underset{\mu}{\text{ARG MIN}} C(\hat{\lambda} + \hat{\lambda}_\mu)$

have simple analytic solutions, because of the special form of C_{Ada} .

→ Note: Statisticians and Machine Learning researchers are different social communities.

(6) Comparison.

MLE:
log-likelihood

$$\min_{\lambda} \left\{ \sum_i \log (e^{\lambda \cdot \phi(d_i)} + e^{-\lambda \cdot \phi(d_i)}) \right.$$

AdaBoost:

$$\min_{\lambda} \left\{ \sum_i e^{-\omega_i \lambda \cdot \phi(d_i)} \right\}$$

Max-Margin

$$\min_{\lambda} \left\{ \frac{1}{2} \|\lambda\|^2 + C \sum_i \max\{0, 1 - \omega_i \lambda \cdot \phi(d_i)\} \right\}$$

Also motivated by
margin ideas.

Also, like AdaBoost, designed to enable
efficient learning algorithm.

Optimization Techniques:

(1) steepest Descent / Discrete Iterative Algorithms.

(2) Coordinate Descent (e.g. AdaBoost)

(3) Online Learning.

→ select a sample x_i at random,
→ do iteration of steepest descent.

High level Points.

The selection of features in d_i is
for AdaBoost than for learning MRF (Della Pietra et al.,
Zhu et al.)
Because AdaBoost only learn a distribution
on one variable ω , AND AdaBoost uses a
simpler cost function.