# Probability Density Estimation

Machine Learning - 10601

Geoff Gordon, Miroslav Dudík
(partly based on slides of Carlos Guestrin and Tom Mitchell)
http://www.cs.cmu.edu/~ggordon/10601/

September 14, 2009

---

## Last time…

**Inference in factor graphs/Bayes nets:**

P(M,Ra,O,W,Ru) $\propto$ ɸ(M) ɸ(Ra) ɸ(O) ɸ(Ra,O,W) ɸ(M,W,Ru)
1  2  3  4  5  $= /Z$

P(W | Ra=F, Ru=T) = ?

**(1) Incorporate evidence:**

P(M, Ra=F, O, W, Ru=T) $\propto$  . . . .

**(2) Eliminate nuisance nodes:**

P(W, Ra=F, Ru=T) $\propto \sum_M \sum_O \cdots = \phi'(w)$

**(3) Normalize:**

P(W | Ra=F, Ru=T) = $\phi'(w) / \sum_w \phi'(W=w)$

# Last time…

**Benefits of factored representations:**

- efficient inference _(sometimes)_
- fewer parameters to estimate ← _approximate_

Δ LEARNING SIMPLER

_less info needed_

---

# Last time: maximum likelihood

heads w/prob **θ**

**N** tosses
**H** heads , N−H tails

$$p(H \mid N, \theta) = \binom{N}{H} \theta^H (1-\theta)^{N-H}$$

← binomial dist.

$$\max_{\theta} \, p(H \mid N, \theta)$$

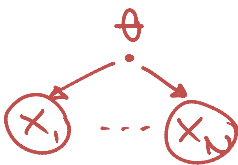$$\hat{\theta}_{ML} = \frac{H}{N}$$

$$\binom{N}{H} = \frac{N!}{H! \, (N-H)!}$$

# Are we learning?

Task: $probability\ estimation\ (prediction)$

Performance measure: $likelihood$

Experience: $thumbtack\ tosses$

$$\theta$$

$$\widehat{x_1} \cdots \widehat{x_N}$$

$$x_n = \begin{cases} 1 & \text{if heads} \quad \text{w/prob } \theta \\ 0 & \text{if tails} \quad \text{w/prob } 1-\theta \end{cases}$$

$$p(x_1, \ldots, x_N | \theta) = \boxed{\prod_n p(x_n | \theta)}$$

$$\theta^{x_n}(1-\theta)^{1-x_n}$$

$$\max_\theta \frac{1}{N} \sum_n \log p(x_n | \theta)$$

$$AVG \approx E_{\theta^*}[\log p(x | \theta)]$$

$$x \sim p(x | \theta^*)$$

$$x_n \sim p(x | \theta^*)$$

$$\text{true}$$

# Are we learning?

$$APPROX$$

$$TRUE\ PERF.\ MEASURE \quad \theta\ fixed$$

$$\frac{1}{N} \sum_n \log p(x_n | \theta) \approx \boxed{E_{\theta^*}[\log p(x | \theta)]}$$

$$x\ random\ var$$

$$\arg\max_\theta E_{\theta^*}[\cdots] = \theta^*$$

# Maximum likelihood estimation

- **expected log likelihood**
  maximized by **true distribution**
- **average log likelihood of data**
  approximates **expected log likelihood**

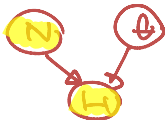~~GREAT!~~ for small sample sizes
approximation poor

# Bayesian approach

initial belief over values of **θ**

e.g. **θ** ∼ uniform over **[0,1]**

$p(\theta) = 1,$ $\qquad \int_0^1 p(\theta)\, d\theta = 1$



$$p(\theta \mid \underline{N}, H) \propto p(\theta, N, H)$$
$$= p(\underline{N})\, p(\theta)\, p(H \mid N, \theta)$$

as a   const.   1   $\binom{N}{H} \theta^H (1-\theta)^{N-H}$
not of θ

$$\propto \theta^H (1-\theta)^{N-H} \quad \text{const.}$$

# Bayesian approach

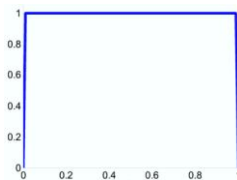**θ, N** parameters

**H** data, $\mathcal{D}$

**p(θ)** prior

**p($\mathcal{D}$ | θ, N)** likelihood

**p(θ | $\mathcal{D}$, N)** posterior

---

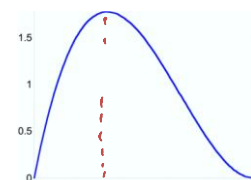# Priors and posteriors

posterior $\propto \theta^H (1-\theta)^{N-H}$

**prior: p(θ)**

**posterior: p(θ | heads=1, tails=1)**

**posterior: p(θ | heads=1, tails=2)**

**posterior: p(θ | heads=20, tails=30)**

MORE OBS
↳ LESS uncertainty
(usually)

# Beta family

$\in [0,1]$

$\mathbf{Beta(\theta \mid \alpha, \beta)} \propto \theta^{\alpha-1}(1-\theta)^{\beta-1}$

- uniform: $\mathbf{p(\theta) = 1}$, $\alpha = \beta = 1$

- after **H** heads, **T** tails: $\mathbf{p(\theta \mid H, T)} \propto \theta^H (1-\theta)^T$

$N-H$, $\alpha = H+1$, $\beta = T+1$
HYPER PARAMS

$p(\mathcal{D} \mid \theta) = \theta^H (1-\theta)^T$

Say $\mathbf{p(\theta) = Beta(\theta \mid \alpha, \beta)} \propto \theta^{\alpha-1}(1-\theta)^{\beta-1}$

- after **H** heads, **T** tails: $\mathbf{p(\theta \mid H, T)}$ $\propto p(\theta, H, T)$

$= p(\theta) p(\mathcal{D} \mid \theta)$

$= \theta^{H+\alpha-1}(1-\theta)^{T+\beta-1}$

$= Beta(H+\alpha, T+\beta)$

# Conjugacy

- if posterior has the same form as prior,
  we say: prior is conjugate relative to likelihood
- e.g.: **Binomial** and **Beta** are conjugate families

Conjugacy: simple Bayesian inference

# Bayesian updating

**p(θ) = Beta(θ | α, β)** *prior*

observe **H** heads, **T** tails

**p(θ | H, T) = Beta(θ | α+H, β+T)** ~~old posterior~~ *new prior*

observe additional **H'** heads, **T'** tails

**p(θ | H, T, H', T')** = $Beta(\theta \mid \underbrace{\alpha+H+H'}_{\alpha'}, \underbrace{\beta+T+T'}_{\beta'})$ *new posterior*
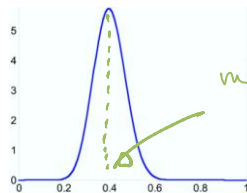
PRIOR = summary of prior experience

---

# Predicting next outcome

observe **H** heads, **T** tails
next observation **X**

- max likelihood $\hat{\theta}_{ML} = \dfrac{H}{H+T}$    $p(X=heads \mid \hat{\theta}_{ML}) = \hat{\theta}_{ML}$

- prior **p(θ) = Beta(θ | α, β)**
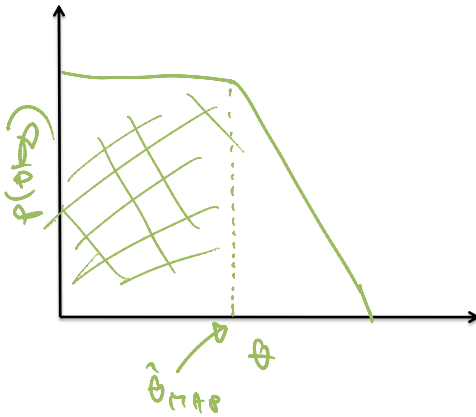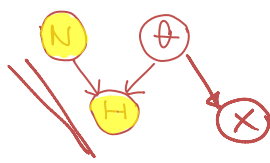  posterior **p(θ | H, T) = Beta(θ | α+H, β+T)**



maximum a posteriori (MAP)
$\hat{\theta}_{MAP}$
if $p(\theta)$ uniform then $\hat{\theta}_{ML} = \hat{\theta}_{MAP}$

# What if posterior looks like…

$p(\theta)$

$\hat{\theta}_{MAP}$

$\theta$

---

# Bayesian prediction

N   θ

H   X

$p(\theta) = \text{Beta}(\theta \mid \alpha, \beta)$

$x = \begin{cases} 1 & \text{if heads} \\ 0 & \text{if tails} \end{cases}$

$P(X \mid N, H)$

$\propto P(x, N, H)$

$= \int P(X, N, \theta, H)\, d\theta$

$= \int \underbrace{P(N, H)}_{\text{const.}} \underbrace{P(\theta \mid N, H)}_{\propto \theta^{H+\alpha-1}(1-\theta)^{T+\beta-1}} \underbrace{P(X \mid N, \theta, H)}_{\propto \theta^{x}(1-\theta)^{1-x}}\, d\theta$

$\propto \int \theta^{H+\alpha+x-1}(1-\theta)^{T+\beta+(1-x)-1}\, d\theta$

# Bayesian prediction

LAPLACE SMOOTHING
(btw UNIFORM & BINOM)

$$P(x|N,H) \propto \int \theta^{H+\alpha+x-1} (1-\theta)^{T+\beta+(1-x)-1} d\theta$$

Beta dist

NORMALIZATION OF BETA

$$\int \theta^{\alpha'-1} (1-\theta)^{\beta'-1} d\theta = \frac{(\alpha'-1)!(\beta'-1)!}{(\alpha'+\beta'-1)!}$$
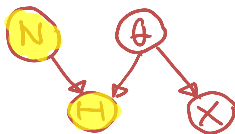
$(\alpha'-1) + (\beta'-1) + 1$

BUT NOT OF
$P(x|N,H)$

$(H+1)H!$

$$P(x|N,H) \propto \begin{cases} \frac{(H+1)! \, T!}{(H+T+2)!} & \text{for } x=1 \\ \frac{H! \, (T+1)!}{(H+T+2)!} & \text{for } x=0 \end{cases}$$

$\frac{(T+1)T!}{}$

$$P(x|b,H) = \begin{cases} \frac{H+1}{H+T+2} & \text{for } x=1 \\ \frac{T+1}{H+T+2} & \text{for } x=0 \end{cases}$$

NORMALIZE → $H+T+2$



# Bayesian prediction vs MAP

**MAP:** $\max_\theta p(\theta|N,H)$

**Bayesian prediction:**

$p(x|\theta, N, H)$   const. $p(U,H)$

$$p(X|N,H) \propto \int p(X|\theta) \, p(\theta|N,H) \, d\theta$$

$$\propto p(x,N,H,\theta)$$

DIFFICULT TO CALCULATE: APPROXIMATE

# What you should know

## Maximum likelihood estimation (MLE)

- approximates maximization of **expected log likelihood**
- **expected log likelihood** maximized by **true distribution**
- approximation poor on small data sets

## Bayesian posterior, MAP, Bayesian prediction

- **posterior** reflects uncertainty in the parameter
- **conjugacy**: posterior has the same form as the prior
  e.g., **Beta** and **Binomial**
- **prior** as a summary of **previous experience (observations)**
- **maximum a posteriori** can suffer similar problems as MLE → STAT ISSUES
- **Bayesian prediction** can be intractable   COMPUT. ISSUES