

Probabilistic Inference Basics

Gaurav Sharma

University of Rochester

Probabilistic Inference Problems

- ▶ We would like to make the "best" inference or "best" estimate of parameters
 - ▶ given observed data x and any prior knowledge we have
 - ▶ typically under assumed model
- ▶ What's best?
 - ▶ Inference: minimize probability of error E

$$\Pr(E) \tag{1}$$

- ▶ Estimation: most likely values of model parameters θ
 - ▶ Alternative formulations as minimization of an error metric

Optimal Decision

- ▶ Min $\Pr(E)$
- ▶ Possible decisions $1, 2, \dots, K$
- ▶ View from perspective of space of possible observations x
 - ▶ Need to partition space of possible observations into K decision regions Z_1, Z_2, \dots, Z_K
 - ▶ Decision i in region Z_i
 - ▶ Problem reduces to partitioning of possible space of observations
 - ▶ What areas should correspond to Z_i where decision is i ?

MAP Decision Optimality

- ▶ Min $\Pr(E)$ equivalent to Max $\Pr(C)$ (Correct)

$$\begin{aligned}\Pr(C) &= \sum_{i=1}^K \Pr(C, i) \\ &= \sum_{i=1}^K \Pr(C | i) p(i) \\ &= \sum_{i=1}^K \int_{Z_i} p(x | i) p(i) dx\end{aligned}$$

- ▶ Optimal rule: For given observation x choose decision that maximizes the argument of integral

$$\hat{i} = \arg \max_i p(x | i) p(i)$$

MAP Decision Optimality

- ▶ MAP Nomenclature
- ▶ Rule for Min $\Pr(E)$

$$\hat{i} = \arg \max_i p(x | i) p(i) \equiv \arg \max_i p(i | x)$$

- ▶ A posteriori probability $p(i | x)$
- ▶ Maximum a posteriori probability (MAP) rule
 - ▶ MAP decisions minimize probability of error
- ▶ Intuitive: choose most likely decision given the data
- ▶ Computationally often challenging: likelihood and prior, prior often unknown
 - ▶ Sometimes likelihood is also challenging to formulate
- ▶ Why likelihood and prior partitioning?

MAP Estimation

- ▶ Model for data implicit in MAP decision
 - ▶ Including model parameters θ
- ▶ Often parameters are unknown a priori, need to also be estimated
- ▶ MAP estimates of parameters

$$\hat{\theta} = \arg \max_{\theta} p(\theta | x) \equiv \arg \max_{\theta} p(x | \theta) p(\theta)$$

- ▶ Intuitive: choose most likely value of parameters given the data
- ▶ Computationally often challenging: likelihood and prior, prior often unknown

Maximum Likelihood (ML) Decision and Parameter Estimation

- ▶ Prior often unknown in decision and estimation problems
 - ▶ Common assumption: equiprobable prior
- ▶ Recall: Posterior probability = likelihood \times prior
- ▶ Under equiprobable prior: maximizing posterior probability \equiv maximizing likelihood
 - ▶ Maximum likelihood (ML) decision/parameter estimation

MAP vs ML Decision

- ▶ MAP Decision

$$\hat{i} = \arg \max_i p(x | i) p(i) \equiv \arg \max_i p(i | x)$$

- ▶ ML Decision

$$\hat{i} = \arg \max_i p(x | i)$$

- ▶ Likelihood often available from "forward" model
- ▶ MAP and ML decisions coincide for equiprobable priors

MAP vs ML Parameter Estimation

- ▶ MAP Estimate

$$\hat{\theta} = \arg \max_{\theta} p(\theta | x) \equiv \arg \max_{\theta} p(x | \theta) p(\theta)$$

- ▶ ML Estimate

$$\hat{\theta} = \arg \max_{\theta} p(x | \theta)$$

- ▶ Likelihood often available from "forward" model
- ▶ MAP and ML estimates coincide for equiprobable priors

ML Parameter Estimation: Bernoulli Example

- ▶ X_i iid Bernoulli, with unknown parameter $\theta = \Pr\{X_i = 1\}$
- ▶ Observations $\mathbf{x} = [x_1, x_2, \dots, x_N]$ string of 0/1 values, length N
- ▶ ML Estimate of θ ?
- ▶ Likelihood function $p(\mathbf{x} \mid \boldsymbol{\theta})$?
 - ▶ Example: $\mathbf{x} = 1101001001$, what is $p(\mathbf{x} \mid \boldsymbol{\theta})$?

ML Parameter Estimation: Bernoulli Example

- Likelihood function

$$p(\mathbf{x} \mid \theta) = \theta^{\sum_i x_i} (1 - \theta)^{N - \sum_i x_i} = \theta^{t(\mathbf{x})} (1 - \theta)^{N - t(\mathbf{x})}$$

- Note $t(\mathbf{x}) = \sum_i x_i$ is a sufficient statistic
 - θ is conditionally independent of \mathbf{x} given $t(\mathbf{x})$
- ML Estimate
 - Calculus of variations

$$t(\mathbf{x})\theta^{t(\mathbf{x})-1}(1 - \theta)^{N-t(\mathbf{x})} - (N - t(\mathbf{x}))\theta^{t(\mathbf{x})}(1 - \theta)^{N-t(\mathbf{x})-1} = 0$$

$$t(\mathbf{x})(1 - \theta) - (N - t(\mathbf{x}))\theta = 0$$

$$\hat{\theta} = \frac{t(\mathbf{x})}{N}$$

- Intuitively appealing: probability of heads is estimated as the empirical fraction of observed heads

ML Parameter Estimation: Scalar Gaussian Example

- ▶ X_i iid Gaussian, with unknown mean μ and variance σ^2
- ▶ Parameters $\theta = [\mu, \sigma^2]$
- ▶ Observations $\mathbf{x} = [x_1, x_2, \dots, x_N]$, sequence of real values
- ▶ ML Estimate of θ ?
- ▶ Likelihood function $p(\mathbf{x} | \theta)$?

$$p(\mathbf{x} | \theta) = \prod_{i=1}^N p(x_i | \theta) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

$$\begin{aligned}\ln p(\mathbf{x} | \theta) &= \sum_{i=1}^N \left(-\frac{1}{2} \ln(2\pi\sigma^2) - \frac{(x_i - \mu)^2}{2\sigma^2} \right) \\ &= -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{\sum_{i=1}^N (x_i - \mu)^2}{2\sigma^2}\end{aligned}$$

ML Parameter Estimation: Scalar Gaussian Example

- ▶ Log Likelihood function

$$\ln p(x | \theta) = -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{\sum_{i=1}^N (x_i - \mu)^2}{2\sigma^2}$$

- ▶ ML Estimates: $\hat{\theta} = [\hat{\mu}, \hat{\sigma}^2] = \arg \max_{\theta} \ln p(x | \theta)$
 - ▶ Calculus of variations

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

- ▶ Intuitively appealing ML estimates are sample mean and sample variance

ML Parameter Estimation: Multivariate Gaussian

- ▶ Multivariate Gaussian X_i , d -dimensional iid Gaussian random vectors
 - ▶ with unknown parameters $\theta = [\mu, \Sigma]$, $d \times 1$ mean vector μ and $d \times d$ covariance matrix Σ
- ▶ Observations $x = [x_1, x_2, \dots, x_N]$, sequence of N , $d \times 1$ vectors
- ▶ Likelihood function

$$\begin{aligned} p(x | \theta) &= \prod_{i=1}^N \mathcal{N}(x_i | \mu, \Sigma) \\ &= \prod_{i=1}^N \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp \left(-\frac{(x_i - \mu)^T \Sigma^{-1} (x_i - \mu)}{2} \right) \end{aligned}$$

$$\ln p(x | \theta) = \sum_{i=1}^N \ln \mathcal{N}(x_i | \mu, \Sigma)$$

ML Parameter Estimation: Multivariate Gaussian

► Log Likelihood function

$$\begin{aligned}\ln p(\mathbf{x} \mid \boldsymbol{\theta}) &= \sum_{i=1}^N \ln \mathcal{N}(x_i \mid \mu, \Sigma) \\ &= \sum_{i=1}^N \ln (\mathcal{N}(x_i \mid \mu, \Sigma)) \\ &= -\frac{N}{2} \ln \left((2\pi)^d \det(\Sigma) \right) \\ &\quad - \frac{1}{2} \sum_{i=1}^N \left((x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right)\end{aligned}$$

ML Parameter Estimation: Multivariate Gaussian

- ▶ Maximization of log likelihood
 - ▶ Tedious but straightforward
 - ▶ Matrix derivatives notation helps see Matrix Cookbook online
- ▶ ML Estimates
 - ▶ (Joint) solution of optimality conditions

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$$
$$\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})(x_i - \hat{\mu})^T$$

- ▶ Intuitively appealing ML estimates are sample mean (vector) and sample covariance (matrix)

ML Parameter Estimation: Homogeneous Markov Chain

- ▶ X_i Homogeneous Markov Chain, with unknown parameters
 $p^1 = [\Pr\{X_1 = i\}]$, $P_{ij} = \Pr\{X_{n+1} = j \mid X_n = i\}$
- ▶ Observations $x = [x_1, x_2, \dots, x_N]$ string of values $\in \{1, 2, \dots, L\}$, length N
- ▶ ML Estimate of θ ?
- ▶ Likelihood function $p(x \mid \theta)$?

ML Parameter Estimation: Homogeneous Markov Chain: Example $L = 2$

- ▶ X_i Homogeneous Markov Chain, with unknown parameters
 $p^1 = [\Pr\{X_1 = i\}]$, $P_{ij} = \Pr\{X_{n+1} = j \mid X_n = i\}$
 - ▶ P is defined by α, β
 - ▶ Depending on assumptions also p^1
- ▶ Observations $x = [x_1, x_2, \dots, x_N]$ string of 0/1 values, length N

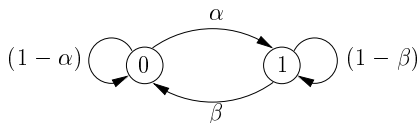


Figure: Transition Diagram representation of a Binary Markov Chain.

- ▶ ML Estimate of θ ?
- ▶ Likelihood function $p(x \mid \theta)$?
 - ▶ Example: $x = 1101001001$, what is $p(x \mid \theta)$?

ML Parameter Estimation: Homogeneous Markov Chain: Example $L = 2$

- ▶ Likelihood function $p(x | \theta)$
 - ▶ Example: $x = 1101001001$, what is $p(x | \theta)$?

$$p(x | \theta) = p^1(1)(1 - \beta)\beta\alpha\beta(1 - \alpha)\alpha\beta(1 - \alpha)\alpha \quad (2)$$

$$= p^1(1)\alpha^{T_{0 \rightarrow 1}(x)}(1 - \alpha)^{(N_0(x) - T_{0 \rightarrow 1}(x))} \times \\ \beta^{T_{1 \rightarrow 0}(x)}(1 - \beta)^{(N_1(x) - T_{1 \rightarrow 0}(x))} \quad (3)$$

- ▶ ML Estimates for α and β : Analogous to Bernoulli case

$$\hat{\alpha} = \frac{T_{0 \rightarrow 1}(x)}{N_0(x)} \quad (4)$$

$$\hat{\beta} = \frac{T_{1 \rightarrow 0}(x)}{N_1(x)} \quad (5)$$

- ▶ What about estimate of p^1 ?

ML Parameter Estimation: General Homogeneous Markov Chain

- ▶ X_i Homogeneous Markov Chain, with unknown parameters
 $p^1 = [\Pr\{X_1 = i\}]$, $P_{ij} = \Pr\{X_{n+1} = j \mid X_n = i\}$
- ▶ Observations $x = [x_1, x_2, \dots, x_N]$ string of values $\in \{1, 2, \dots, L\}$, length N
- ▶ ML Estimate of transition probability matrix P :
 - ▶ Calculus of variations (with constraints)

$$\hat{P}_{ij} = \frac{T_{i \rightarrow j}(x)}{N_i(x)} \quad (6)$$

$T_{i \rightarrow j}(x) = \#$ of $i \rightarrow j$ transitions in x

$N_i(x) = \#$ of transitions in x originating from state i

Outlook

- ▶ Using models for probabilistic inference and estimation
- ▶ Expectation Maximization
 - ▶ Modeling interactions that are not directly visible
 - ▶ Hidden/latent variables
 - ▶ Two case studies
 - ▶ WAMI to Roadmap alignment + Flow cytometry cell clustering
 - ▶ IID latent variables
 - ▶ General EM formulation and Gaussian mixture model

Outlook II

- ▶ Hidden Markov Models (HMMs)
 - ▶ Build upon Markov models
 - ▶ Memory + hidden/latent variables
 - ▶ Two case studies
 - ▶ Sequence alignment + Error correction decoding for convolutional codes
 - ▶ Three standard problems
 - ▶ Sequence estimation, individual state marginal probability estimation, parameter estimation
 - ▶ Parameter estimation: Baum-Welch \equiv EM

Outlook III

- ▶ Stochastic context free grammars
 - ▶ Generalization of dependency beyond HMM Markovian structure
 - ▶ Example: Palindromes. Are they modeled well by a HMM?
 - ▶ Three standard problems analogous to HMMs
 - ▶ Sequence estimation (CYK), individual state marginal probability estimation, parameter estimation
 - ▶ Parameter estimation: Inside-Out \equiv EM

Outlook IV

- ▶ Markov Random Fields
 - ▶ Generalization of dependency more appropriate for multi-dimensional data
 - ▶ Graphical representation
- ▶ Additional relevant examples from current ongoing research
 - ▶ Turbo decoding in communications and RNA structure prediction
 - ▶ Color barcodes