

- Previously assumed we knew priors  $p(y)$  and  $f_{xy}(x|y)$
- But rarely have complete knowledge
- May have general knowledge plus samples / training data
- Use to design / train classifier.
- Use samples to estimate unknown probabilities / pdfs and use estimates as true values.
- Estimation of  $f_{xy}(x|y)$  problematic when # samples too small and when dimensionality of  $x$  is large
- If can parameterize conditional PDF, e.g.  
assume  $f_{xy}(x|i) \sim \mathcal{N}(\mu_i, \Sigma_i)$  (without knowing  $\mu_i, \Sigma_i$ ) simplifies problem to parameter estimation
- Maximum Likelihood Estimation - parameters with fixed unknown values: best estimate is one which maximizes prob. of obtaining observed samples.
- Bayesian methods: parameters as r.v.'s having known prior distribution.  
Observation of samples converts this to a posterior density. Typically, additional samples sharpen a posterior density, peak near true values

## ML estimation

- have good convergence properties as # training samples increases
- simpler than alternative methods
- use set  $D$  of training samples drawn indep. from  $f(x|\underline{\theta})$  to estimate unknown parameter  $\underline{\theta}$   
e.g.  $\underline{\theta} = (\mu, \Sigma)$
- Suppose  $D$  contains samples  $x_1, \dots, x_n$

$$f(D|\underline{\theta}) = \prod_{k=1}^n f(x_k|\underline{\theta})$$

likelihood of  $\underline{\theta}$  wrt ~~same~~ set of samples

- ML estimate of  $\underline{\theta}$  is  $\hat{\underline{\theta}}_{ML}$

$$\hat{\underline{\theta}}_{ML} = \underset{\underline{\theta}}{\operatorname{argmax}} f(D|\underline{\theta})$$

value of  $\underline{\theta}$  that best agrees/supports actually observed training samples.

- $\ln(\cdot) \rightarrow$  usu. easier to work with

log likelihood function

$$\begin{aligned}\hat{\underline{\theta}}_{ML} &= \underset{\underline{\theta}}{\operatorname{argmax}} \ln f(D|\underline{\theta}) \\ &= \underset{\underline{\theta}}{\operatorname{argmax}} \ln \prod_k f(x_k|\underline{\theta}) \\ &= \underset{\underline{\theta}}{\operatorname{argmax}} \sum_{k=1}^n \ln f(x_k|\underline{\theta})\end{aligned}$$

Ex: Data: observed seq.  $D$  of  $n_H$  heads and  $n_T$  tails

Model: Each flip follows Bernoulli distribution

$$P(H) = \theta, \quad P(T) = 1 - \theta, \quad \theta \in [0, 1].$$

likelihood of observing seq.  $D$  is

$$P(D|\theta) = \theta^{n_H} (1-\theta)^{n_T}$$

Given model and data, estimate  $\theta$ .

ML estimate :  $\hat{\theta}_{ML} = \underset{\theta}{\operatorname{argmax}} P(D|\theta)$

$$\begin{aligned} &= \underset{\theta}{\operatorname{argmax}} \ln P(D|\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \ln (\theta^{n_H} (1-\theta)^{n_T}) \\ &= \underset{\theta}{\operatorname{argmax}} \underbrace{n_H \ln \theta + n_T \ln(1-\theta)}_{\text{concave fn. of } \theta} \end{aligned}$$

$$\frac{d}{d\theta} \ln P(D|\theta) = 0:$$

$$\frac{n_H}{\theta} - \frac{n_T}{1-\theta} = 0$$

$$\Rightarrow \hat{\theta}_{ML} = \frac{n_H}{n_H + n_T}$$

So e.g. flip 10 times, get H 8 times, T 2 times

$$\text{Then } \hat{\theta}_{ML} = \frac{n_H}{n_H + n_T} = 0.8$$

Trusting data completely. There could be too little or noisy data.

Ex: Data: observe  $D = \{x_1, \dots, x_n\}$

Model:  $x_i \sim \mathcal{N}(\mu, \sigma^2)$  iid

Estimate:  $\mu, \sigma^2$

$$\hat{\theta}_{ML} = \underset{\theta}{\operatorname{argmax}} \ln f(D|\theta)$$

$$= \text{" } \ln \prod_{k=1}^n f(x_k|\theta)$$

$$= \text{" } \ln \prod_{k=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x_k - \mu)^2}{2\sigma^2}\right]$$

$$= \text{" } \sum_{k=1}^n -\frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln \sigma^2 - \frac{(x_k - \mu)^2}{2\sigma^2}$$

$$= \underset{\theta}{\operatorname{argmin}} \sum_{k=1}^n \frac{1}{2} \ln(2\pi) + \frac{1}{2} \ln \sigma^2 + \frac{(x_k - \mu)^2}{2\sigma^2}$$

convex in  $\mu$

$$\frac{\partial \ln f(D|\theta)}{\partial \mu} = \sum_{k=1}^n \frac{(x_k - \mu)}{\sigma^2} = 0$$

$$\Rightarrow \hat{\mu}_{ML} = \frac{1}{n} \sum_{k=1}^n x_k \quad \text{sample mean}$$

$$\frac{\partial \ln f(D|\theta)}{\partial \sigma^2} = 0: \quad \hat{\sigma}_{ML}^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu}_{ML})^2$$

Is  $\hat{\mu}_{ML}$  unbiased?

$E[\hat{\mu}_{ML}] \stackrel{?}{=} \mu$  i.e. expected value of estimator over all datasets of size  $n$  = true value

$$E[\hat{\mu}_{ML}] = E\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n} \sum_{i=1}^n E[x_i] = \frac{1}{n} \sum_{i=1}^n \mu = \mu \quad \checkmark$$

unbiased

Is  $\hat{\sigma}_{ML}^2$  unbiased?

$$E[\hat{\sigma}_{ML}^2] = E\left[\frac{1}{n} \sum_{k=1}^n (X_k - \hat{\mu}_{ML})^2\right]$$

$$= E\left[\frac{1}{n} \sum_{k=1}^n ((X_k - \mu) + (\mu - \hat{\mu}))^2\right]$$

$$= E\left[\frac{1}{n} \sum_{k=1}^n \{(X_k - \mu)^2 + 2(X_k - \mu)(\hat{\mu} - \mu) + (\hat{\mu} - \mu)^2\}\right]$$

$$= E\left[\frac{1}{n} \sum_{k=1}^n (X_k - \mu)^2\right] - 2E\left[\frac{1}{n} \sum_{k=1}^n (X_k - \mu)(\hat{\mu} - \mu)\right] + E\left[(\hat{\mu} - \mu)^2\right]$$

$$= \frac{1}{n} \sum_{k=1}^n \sigma^2 - 2E\left[\left(\frac{1}{n} \sum_{k=1}^n X_k - \frac{1}{n} \sum_{k=1}^n \mu\right)(\hat{\mu} - \mu)\right] + E\left[(\hat{\mu} - \mu)^2\right]$$

$$= \sigma^2$$

$$- 2E[(\hat{\mu} - \mu)(\hat{\mu} - \mu)]$$

$$- 2E[(\hat{\mu} - \mu)^2]$$

$$- E[(\hat{\mu} - \mu)^2]$$

$$E[\hat{\sigma}_{ML}^2] = \sigma^2 - \frac{\sigma^2}{n}$$

$$= \frac{(n-1)\sigma^2}{n}$$

biased

$$\text{var}\left(\frac{1}{n} \sum_{k=1}^n X_k\right)$$

$$= \frac{1}{n^2} n \sigma^2 = \frac{\sigma^2}{n}$$

However  $\lim_{n \rightarrow \infty} E[\hat{\sigma}_{ML}^2] = \sigma^2$

asymptotically unbiased : estimator becomes unbiased in limit as # samples  $\rightarrow \infty$

## Properties of ML estimators

- Assume data generated from some true parameter value  $\underline{x} \sim f(\underline{x}|\theta^*)$ . ML estimator  $\hat{\theta}_{ML}$  obtained from data set of size  $n$  satisfies.
  - asymptotically unbiased  $\lim_{n \rightarrow \infty} E[\hat{\theta}_{ML}] = \theta^*$  ~~as  $n \rightarrow \infty$~~
  - Variance ( $\hat{\theta}_{ML}$ ) is minimum among estimators
  - If  $\hat{\theta}_{ML}$  is for  $\theta$ , then  $g(\hat{\theta}_{ML})$  is ML estimator for  $g(\theta)$ .
  - Distribution of  $\hat{\theta}_{ML}$  is Gaussian for large  $n$ .

## Multivariate Gaussian

$$D = (\underline{x}_1, \dots, \underline{x}_n)$$

$$f(D|\mu, \Sigma) = \prod_{k=1}^n f(\underline{x}_k|\mu, \Sigma)$$

$$= \prod_{k=1}^n \frac{1}{(2\pi)^{d/2} \sqrt{\det(\Sigma)}} \exp\left[-\frac{1}{2} \frac{(\underline{x}_k - \mu)^T \Sigma^{-1} (\underline{x}_k - \mu)}{1}\right]$$

$$\ln f(D|\mu, \Sigma) = \sum_{k=1}^n -\frac{1}{2} \ln[(2\pi)^d \det(\Sigma)]$$

$$- \frac{1}{2} (\underline{x}_k - \mu)^T \Sigma^{-1} (\underline{x}_k - \mu)$$

$$\frac{\partial}{\partial \mu} = 0$$

$$\hat{\mu}_{ML} = \frac{1}{n} \sum_{k=1}^n \underline{x}_k$$

$$\frac{\partial}{\partial \Sigma} = 0:$$

$$\hat{\Sigma}_{ML} = \frac{1}{n} \sum_{k=1}^n (\underline{x}_k - \hat{\mu}_{ML})(\underline{x}_k - \hat{\mu}_{ML})^T$$

## MAP Estimation

-  $\underline{\theta}$  considered as a random vector; Have prior belief

$$\hat{\underline{\theta}}_{\text{MAP}} = \underset{\underline{\theta}}{\operatorname{argmax}} p(\underline{\theta} | D) \quad \begin{array}{l} \text{regarding } p(\underline{\theta}) \\ \text{incorporate data} \end{array}$$

$$= \underset{\underline{\theta}}{\operatorname{argmax}} \frac{p(D | \underline{\theta}) p(\underline{\theta})}{p(D)} \quad D = \{x_1, \dots, x_n\}$$

$$= \underset{\underline{\theta}}{\operatorname{argmax}} \underbrace{p(D | \underline{\theta})}_{\text{likelihood}} \underbrace{p(\underline{\theta})}_{\text{prior}}$$

$$= \underset{\underline{\theta}}{\operatorname{argmax}} \left( \prod_{k=1}^n p(x_k | \underline{\theta}) \right) p(\underline{\theta})$$

$$= \underset{\underline{\theta}}{\operatorname{argmax}} \sum_{k=1}^n \left( \ln p(x_k | \underline{\theta}) + \ln p(\underline{\theta}) \right)$$

Ex: Assume  $f(x | \mu) = \mathcal{N}(x, \sigma^2)$  where  $\mu$  is a r.v. with prior  $f(\mu)$   
 $\uparrow$   
known

$$\underbrace{f(\mu)}_{\text{prior}} = \mathcal{N}(\underbrace{\mu_0}_{\uparrow \text{known}}, \underbrace{\sigma_0^2}_{\uparrow})$$

$\mu_0$ : best prior guess for  $\mu$

$\sigma_0^2$ : uncertainty in prior guess for  $\mu$

$D = (x_1, \dots, x_n)$  are observed samples

$$\hat{\mu}_{\text{MAP}} = \underset{\mu}{\operatorname{argmax}} \underbrace{f(\mu | D)}_{\text{posterior}}$$

$$= \operatorname{argmax}_{\mu} f(\mu) f(\mu)$$

$$= \prod_{k=1}^n f(x_k | \mu) f(\mu)$$

$$= \prod_{k=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_k - \mu)^2}{2\sigma^2}} \cdot \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{(\mu - \mu_0)^2}{2\sigma_0^2}}$$

$$= -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{k=1}^n (x_k - \mu)^2$$

$$- \frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln \sigma_0^2 - \frac{1}{2\sigma_0^2} (\mu - \mu_0)^2$$

$$\frac{\partial}{\partial \mu} = 0: \quad \frac{1}{\sigma^2} \sum_{k=1}^n (x_k - \mu) = \frac{1}{\sigma_0^2} (\mu - \mu_0)$$

$$\hat{\mu}_{MAP} = \frac{\sigma_0^2}{\sigma^2 + \sigma_0^2 n} \sum_{k=1}^n x_k + \frac{\sigma^2}{\sigma^2 + \sigma_0^2 n} \mu_0$$

$$= \frac{\sigma_0^2 n}{\sigma^2 + \sigma_0^2 n} \left( \underbrace{\frac{1}{n} \sum_{k=1}^n x_k}_{\hat{\mu}_{ML} \text{ data}} \right) + \underbrace{\frac{\sigma^2}{\sigma^2 + \sigma_0^2 n} \mu_0}_{\text{prior}}$$

Note  $\hat{\mu}_{MAP}$  bet.  $\mu_0$  and  $\hat{\mu}_{ML}$

$\uparrow$                        $\uparrow$   
 prior                      data

• If  $\sigma_0 \neq 0$ , as  $n \rightarrow \infty$ ,  $\hat{\mu}_{MAP} \rightarrow \hat{\mu}_{ML}$

• If  $\sigma_0 = 0$ ,  $\hat{\mu}_{MAP} = \mu_0$

- degenerate situation where prior is certain



- If  $\sigma_0^2 \gg \sigma$ ,  $\hat{\mu}_{MAP} \sim \hat{\mu}_{ML}$  : high uncertainty about prior
- If  $\sigma^2/\sigma_0^2$  is finite, after enough samples  $\hat{\mu}_{MAP}$  will converge to  $\hat{\mu}_{ML}$  regardless of  $\mu_0$  and  $\sigma_0^2$ .

Ex: coin flip example:  $D = \{ \alpha, H, \alpha_0 T's \}$

Suppose have prior  $f(\theta) = \text{Beta}(\beta_0, \beta_1) = \frac{\theta^{\beta_1-1} (1-\theta)^{\beta_0-1}}{B(\beta_0, \beta_1)}$

$\beta_1, \beta_0$  indicate fractions and H and T . normalization

$$\begin{aligned} \hat{\theta}_{MAP} &= \underset{\theta}{\text{argmax}} \quad P(D|\theta) f(\theta) \\ &= \text{"} \quad \theta^{\alpha_1} (1-\theta)^{\alpha_0} \frac{\theta^{\beta_1-1} (1-\theta)^{\beta_0-1}}{B(\beta_0, \beta_1)} \\ &= \text{"} \quad \frac{\theta^{\alpha_1+\beta_1-1} (1-\theta)^{\alpha_0+\beta_0-1}}{B(\beta_0, \beta_1)} \\ &= \text{"} \quad \theta^{\alpha_1+\beta_1-1} (1-\theta)^{\alpha_0+\beta_0-1} \end{aligned}$$

$$= \frac{(\alpha_1 + \beta_1 - 1)}{(\alpha_1 + \beta_1 - 1) + (\alpha_0 + \beta_0 - 1)}$$

Add  $\beta_1 - 1$   
imag. heads  
 $\beta_0 - 1$  imag.  
tails

$$\text{vs. } \hat{\theta}_{ML} = \frac{\alpha_1}{\alpha_1 + \alpha_0}$$

why choose beta prior, b/c a posterior is also beta  
conjugate prior