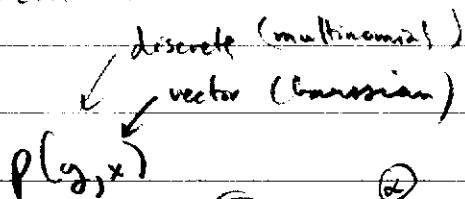


-classifier



i.e. for M classes have
 M Gaussians

$$p(x, y | \theta) = \left(\prod_j \alpha_j^{y_j} \right) \left(\sum_m \delta(y=m) N(x | \mu_m, \Sigma_m) \right)$$

$$\sum \log p(x_i, y_i | \theta) = \sum \log p(y_i, x_i | \theta)$$

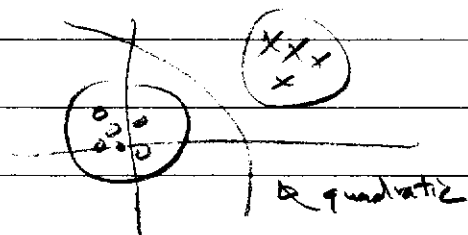
$$= \sum_i \log \left(\prod_j \alpha_j^{y_{ij}} \sum_m \delta(m=y_{ij}) N(x_i | \mu_m, \Sigma_m) \right)$$

$$= \sum_i \sum_j y_{ij} \log \alpha_j + \sum_i \log \sum_m \delta(m=y_{ij}) N(x_i | \mu_m, \Sigma_m)$$

$$\frac{2}{2\mu_m}$$

$$\sum \delta(m=y_{ij}) \log N(x_i | \mu_m, \Sigma_m)$$

$$\mu_m = \frac{\sum x_i \delta(m=y_{ij})}{\sum \delta(m=y_{ij})}$$

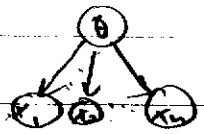


- Readings: Chapter 2 (first half)

Chapter 4

Chapter 7 (next week)

- Quick Review of Graphical Models: $p(x_i | p(x_i))$



- Quick Rev of why pdf (all vars)

- Quick Rev of Bayesian Learning

post = $\frac{\text{likel} \times \text{prior}}{\text{evidence}}$

$p(x|X) \propto \prod p(x_i|\theta) p(\theta) d\theta$



- Quick Rev of ML \rightarrow intuitive mean & cov \rightarrow ML

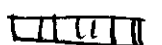
- Questions / Conversations / Other Distributions \Rightarrow e-family

- scalar
- coins
- discrete
- integer

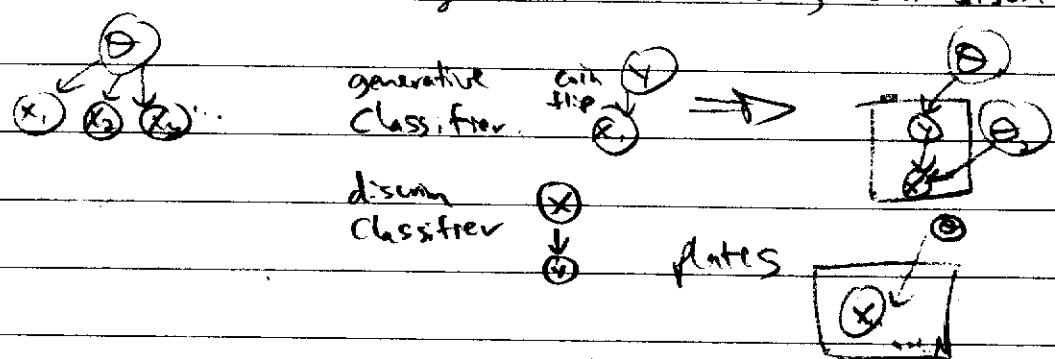
nice properties for ML / Bayes

Review Optimization

- MAP Gaussian & Bayesian Gaussian (page 8-8bis)

- Bernoulli / Multinomial / simplest hypercube  / heuristics then ML derivation (page 7-8)

- Generative Models \rightarrow regenerate the data, later Discrim.



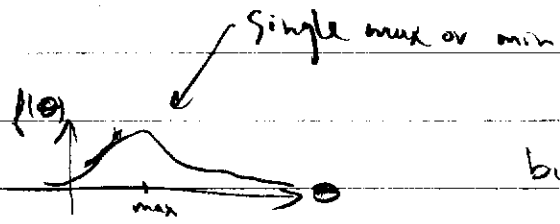
- Classifier $p(y_i|x)$ (page 9)

- When is ML consistent (analytic?) & Bayes? E-family,

E-family: Bernoulli, multinomial, Gaussian, Poisson, Dirichlet,
 binom \downarrow discrete \mathbb{R}^k \mathbb{Z}^+ the \mathbb{R}^n

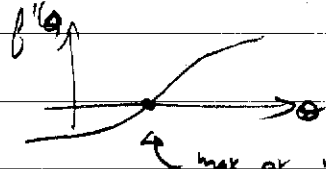
(# of people in a room)

Optimization:



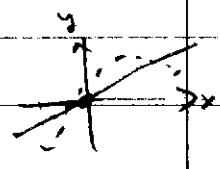
but min if function is a bowl

$$\frac{\partial f}{\partial \theta}$$

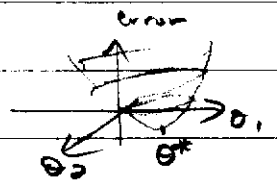


$$\frac{\partial f}{\partial \theta} = 0 \text{ \& solve.}$$

∇f direction of fastest increase



least squares: $\min_{\theta} \frac{1}{2} \sum_i (y_i - \theta^T x_i)^2$



$$\frac{\partial}{\partial \theta} \left[\frac{1}{2} (\vec{y} - X\theta)^T (\vec{y} - X\theta) \right] = 0$$

$$X^T (\vec{y} - X\theta) = 0$$

$$X^T X \theta = X^T \vec{y}$$

$$\theta = (X^T X)^{-1} X^T \vec{y}$$

easy to solve

weighted least squares: $\min_{\theta} \frac{1}{2} \sum_i w_i (y_i - \theta^T x_i)^2$

$$\frac{\partial}{\partial \theta} \left[\frac{1}{2} (\vec{y} - X\theta)^T W (\vec{y} - X\theta) \right] = 0$$

$$X^T W (\vec{y} - X\theta) = 0$$

Hessian gradient

$$W^T X \theta = X^T W \vec{y}$$

$$\theta = (W^T X)^{-1} X^T W \vec{y}$$

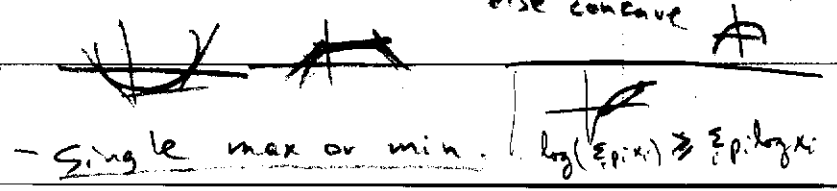
not easy to solve

gradient desc: $\theta^{t+1} = \theta^t + \delta \nabla_{\theta} L$ ← function to maximize

newton-raphson: (Taylor-series) $\theta^{t+1} = \theta^t + H^{-1} \nabla_{\theta} L \Big|_{\theta^t}$

convex/concave: $-\frac{\partial}{\partial \theta} \left(\frac{\partial}{\partial \theta} J(\theta) \right) = H$ if the semi-def then convex else concave

epigraph
- Jensen



natural params...

- General form: $p(x|\theta) = h(x) \exp(\theta^T T(x) - A(\theta))$

$$p(x|\theta) = \exp(H(x) + \theta^T x - A(\theta))$$

$\uparrow H(x) = \log h(x)$

$\uparrow A(\theta)$ convex

$$\int p(x|\theta) dx = 1$$

Cumulant gen. func.

$$A(\theta) = \log \int \exp(H(x) + \theta^T x) dx$$

Laplace transform

$$\frac{\partial}{\partial \theta} A(\theta) = \text{mean} = E[T(x)] := \mu$$

$$\frac{\partial^2}{\partial \theta^2} A(\theta) = \text{tr. def.}, \text{ cov.} = E[T(x)T(x)^T] - (E[T(x)])^2$$

gradient $\frac{\partial}{\partial \theta} A(\theta)$ is unique (convex)

$$\text{to a } \theta \\ E[T(x)] = \mu \iff \theta$$

- Sufficient stats. $\sum T(x)$ or $E[T(x)]$ is all you need

- ML e-family has this, get $E x, E x^2, \dots E x^p$ stop.
- ML $\prod_{i=1}^N \exp(H(x_i) + \theta^T x_i - A(\theta)) : \frac{\partial}{\partial \theta} A(\theta) = \frac{1}{N} \sum x_i$

Examples Bernoulli: $p(x|\pi) = \pi^x (1-\pi)^{1-x}$ coin
 $\exp(\log(\frac{\pi}{1-\pi}) x + \log(1-\pi))$

$$\theta = \frac{\pi}{1-\pi} \implies \text{or } \pi = \frac{e^\theta}{1+e^\theta}$$

$$T(x) = x$$

$$A(\theta) = -\log(1-\pi) = \log(1+e^\theta)$$

$$h(x) = 1 \implies H(x) = 0$$

example multinomial $p(x|\phi) = \prod_{i=1}^m \phi_i^{x_i}$ $x_i \in [0, 1]$
 $\sum_{i=1}^m x_i = 1$
 $p(x|\theta) = \exp\left(\sum_{i=1}^m x_i \ln \phi_i\right)$
 $= \exp\left(\sum_{i=1}^m x_i \ln \phi_i + (1 - \sum_{i=1}^m x_i) \ln(1 - \sum_{i=1}^m \phi_i)\right)$ $\sum_{i=1}^m \phi_i = 1$
 $= \exp\left(\sum_{i=1}^m \ln\left(\frac{\phi_i}{1 - \sum_{j=1}^m \phi_j}\right) x_i + \ln(1 - \sum_{j=1}^m \phi_j)\right)$
 $\theta_i := \ln\left(\frac{\phi_i}{\phi_m}\right)$ $\phi_i = \frac{e^{\theta_i}}{\sum_{j=1}^m e^{\theta_j}}$ $\phi_m = 1 - \sum_{i=1}^{m-1} \phi_i$
 $A(\theta) = \ln\left(1 + \sum_{i=1}^{m-1} e^{\theta_i}\right)$
 $H(\theta) = 0$

example dirichlet $p(x|\alpha) = \frac{\Gamma(\alpha_1 + \dots + \alpha_k)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_k)} \prod_{i=1}^k \phi_i^{\alpha_i - 1}$

?

example Gaussian $\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$ $\alpha > 1$
 $= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(\frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}x^2 - \frac{1}{2\sigma^2}\mu^2 - \ln\sigma\right)$

multi-dim. $\left[\begin{matrix} \mu/\sigma^2 \\ \ln\sigma \\ \frac{1}{2\sigma^2} \end{matrix} \right]$

$\theta = \begin{bmatrix} \mu/\sigma^2 \\ -1/2\sigma^2 \end{bmatrix}$
 $T(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$

✓ this bottom param is neg.

$A(\theta) = \frac{\mu^2}{2\sigma^2} + \ln\sigma = \frac{-\theta_1^2}{4\theta_2} - \frac{1}{2} \ln(-2\theta_2)$
 $H(x) = -\frac{1}{2} \ln(2\pi)$

example exponential $p(x) = \lambda e^{-\lambda x}$ $x \geq 0$ +ve scalars
 $p(x) = \exp(-\lambda x + \ln \lambda)$ $\theta := -\lambda$
 $H(\theta) = 0$
 $\forall \theta < 0$ $A(\theta) = -\ln(-\theta)$

example poisson $p(x|\lambda) = \frac{\lambda^x e^{-\lambda}}{x!}$
 $= \frac{1}{x!} \exp(x \log \lambda - \lambda)$
 $\theta = \log \lambda$
 $T(x) = x$
 $A(\theta) = \lambda = e^\theta$
 $H(x) = -\log x!$

Conjugate prior
 ↓
 Conjugate prior

more general form: $\exp(H(x) + \phi(\theta)^T T(x) - A(\theta))$
 unique parts
 conjugate: $\exp(J(\eta, \nu) + \phi(\theta)^T T(x) - \eta A(\theta))$
 (Legendre transform) ↑ virtual data n_0 of times.

Gauss	Gauss
Mult.	Dir
Bern	Beta
Gauss	Normal
	Inv. Wish.

MAR

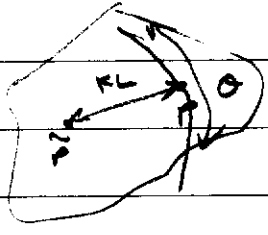
Bayes

$\int p(x|\theta) \prod p(x_i|\theta) p(\theta) d\theta$ is solvable
 ↑ ↑ ↑ conjugate
 e-family

KL-divergence: $\equiv ML$

$\hat{p}(x) = \frac{1}{N} \sum_{i=1}^N \delta(x, x_i)$ $\eta/\eta/\eta$

$\min D(\hat{p} || p) = \sum_x \hat{p}(x) \log \left(\frac{\hat{p}(x)}{p(x|\theta)} \right)$



$= \sum_x \hat{p}(x) \log \hat{p}(x) - \frac{1}{N} \ell(\theta|D)$
 ↑ negative entropy

Assignment on Web Page Tomorrow
 Read: Chapters 9 & 10

- Assignment #1 on the web page, due Tuesday 19th at 5pm
by email or my office, Questions? office hours (email)

- Review: optimization: grad descent, or $\frac{\partial}{\partial \theta} = 0$

e-family: ML maximization step $\frac{\partial}{\partial \theta} = 0$ guaranteed unique
nice form $\prod_{n=1}^N p(x_n | \theta)$ by i.i.d. exp.

e-family: each distrib has a conjugate

conj (conj (e)) = e $\rightarrow p(x|\theta) \rightarrow p(\theta|x)$

Gaussian μ - Gaussian μ

Gaussian σ^2 - Inverse Wishart
Multinomial - Dirichlet

easy Bayes when prior is conjugate

Show examples - multinomial & Gaussian

Then Break e-family \rightarrow EM \rightarrow K-means

today ch. 9 (2) & ch. 10
next ch. 2 (and 4) & ch. 3

1) Conquer & Divide

3) Bound Max

4) Filling in nodes

5) Coordinate Ascent

6) Information Geometry

- Bayesian Gaussian

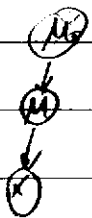
$p(\mu) \rightarrow \mu$

$\mu = \frac{1}{N} \sum_{n=1}^N x_n$ Σ given

see page 8-8bis

changes σ , makes it bigger

changes μ , pulls it towards prior

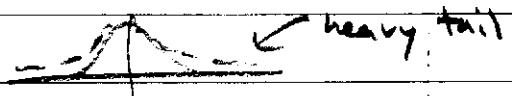


- Full integral over covariance params too

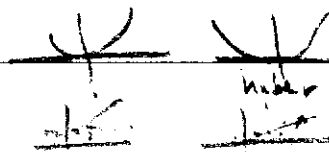
$p(\Sigma) = IW(\Sigma)$

get student "t"

slower than e^{-x^2}

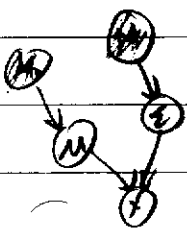


robust stats

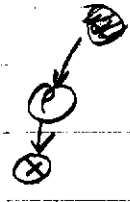


as $N \rightarrow \infty$

Bayes \leftrightarrow ML



Bayesian Multinomial



$$p(x|p) = \prod_{k=1}^K p_k^{\delta(x=k)} \leftarrow \text{indicator function}$$

$$p(p) = \text{Dirichlet} = \frac{\prod_k \Gamma(\alpha_k)}{\Gamma(\sum \alpha_k)} \prod_k p_k^{\alpha_k - 1} \leftarrow \text{hyper params}$$

$\sum_k p_k = 1 \quad p_k \geq 0$

$$\Gamma(n+1) = n!$$

$\approx p(x, X)$

$$p(x|X) = \frac{1}{p(X)} \int p(x|p) \prod_{n=1}^N p(x_n|p) p(p) dp$$

$$= \frac{1}{p(X)} \int \prod_{k=1}^K p_k^{\sum \delta(x_n=k)} \prod_{n=1}^N \prod_{k=1}^K p_k^{\delta(x_n=k)} \frac{\prod_k \Gamma(\alpha_k)}{\Gamma(\sum \alpha_k)} \prod_k p_k^{\alpha_k - 1} dp$$

$$= \frac{1}{p(X)} \frac{\prod_k \Gamma(\alpha_k)}{\prod_k \Gamma(\alpha_k)} \int \prod_k p_k^{\alpha_k - 1 + \sum \delta(x_n=k) + \delta(x=k)} dp$$

$$= \frac{1}{p(X)} \frac{\prod_k \Gamma(\alpha_k)}{\prod_k \Gamma(\alpha_k)} \frac{\Gamma(\sum \alpha_k + \sum \delta(x_n=k) + \delta(x=k))}{\Gamma(\sum \alpha_k + \sum \delta(x_n=k) + \delta(x=k))}$$

easier: $\sum \delta(x_n=k) = N_k$
 $\sum N_k = N$

$$p(X) = \frac{\prod_k \Gamma(\alpha_k)}{\Gamma(\sum \alpha_k)} \frac{\prod_k \Gamma(\alpha_k + N_k)}{\Gamma(N + \sum \alpha_k)}$$

$$p(x|X) = \frac{\prod_k \Gamma(\alpha_k + N_k + \delta(x=k))}{\Gamma(\sum \alpha_k + N + 1)} \frac{\Gamma(N + \sum \alpha_k)}{\prod_k \Gamma(\alpha_k + N_k)}$$

$$p(x=k|X) = \frac{N_k + \alpha_k}{N + \sum \alpha_k}$$

← "plug in" $x=k$
 ← get ... instead of $\frac{N_k}{N}$