

eq. 3.2 $p(y|x, \theta) = \prod_i p(y_i | x_i, \theta)$ where $\theta = \{k, \phi_f\}$
 \hookrightarrow params of form f

Poisson likelihood: $p(x_i) = e^{-\lambda_0} \frac{1}{x_i!} \lambda_0^{x_i}$

let $\lambda_0 \leftarrow \Delta f(k \cdot x_i)$, then:

eq. 3.3 $p(y_i | x_i, \theta) = e^{-\Delta f(k \cdot x_i)} \frac{1}{y_i!} (\Delta f(k \cdot x_i))^{y_i}$

$$p(y|x, \theta) = \prod_i \left[e^{-\Delta f(k \cdot x_i)} \frac{1}{y_i!} (\Delta f(k \cdot x_i))^{y_i} \right]$$

$$= \prod_i \left[e^{-\Delta f(k \cdot x_i)} (\Delta^y f(k \cdot x_i)^y) / y_i! \right]$$

eq. 3.4 $= \Delta^n \prod_i \left[e^{-\Delta f(k \cdot x_i)} f(k \cdot x_i)^y / y_i! \right]$ where $n = \sum_i y_i$

find MLE of $\hat{\theta} = \{\hat{k}, \hat{\phi}_f\}$ by max. log-likelihood:

$$\begin{aligned} \log[p(y|x, \theta)] &= \log \left[\Delta^n \prod_i \left[e^{-\Delta f(k \cdot x_i)} f(k \cdot x_i)^y / y_i! \right] \right] \\ &= \log \Delta^n + \log \left[\prod_i e^{-\Delta f(k \cdot x_i)} \right] + \log \left[\prod_i f(k \cdot x_i)^y \right] - \log \left[\prod_i y_i! \right] \\ &= c_1 + \left[\sum_i \log(e^{-\Delta f(k \cdot x_i)}) + \sum \log(f(k \cdot x_i)^y) \right] + c_2 \end{aligned}$$

eq. 3.5 $= -\Delta \sum_i f(k \cdot x_i) + \sum_i y_i \log(f(k \cdot x_i)) + c$

constrain k to be a unit vector: $\|k\| = 1$

$\Rightarrow k \cdot x_i$ will only change direction of x_i , not magnitude
 now estimate k using "angular error"

$$\frac{d}{dk} \left(\log [p(y|x, \theta)] \right) = \frac{d}{dk} \left[-\Delta \sum_i f(k \cdot x_i) + \sum_i y_i \log(f(k \cdot x_i)) + c \right]$$

$$= -\Delta \sum_i x_i \cdot f'(k \cdot x_i) + \sum_i y_i \cdot x_i \cdot f'(k \cdot x_i) \cdot \frac{1}{f(k \cdot x_i)} + 0$$

$$= -\Delta \sum_i x_i f'(k \cdot x_i) + \sum_i y_i x_i \frac{f'(k \cdot x_i)}{f(k \cdot x_i)}$$

use "method of Lagrange multipliers to find the local max/min:

find max/min of $f(x)$ w/ constraint $g(x) = 0$, then

$$\mathcal{L}(x, \lambda) = f(x) - \lambda g(x)$$

["Lagrange multiplier"]

$$\text{set } \frac{d}{dx} \left[\log(p(y|x, \theta)) \right] = 0 = g(x) \quad (\text{where } x \text{ here is our } k)$$

$$\mathcal{L}(k, \lambda) = \left(-\Delta \sum_i f(k \cdot x_i) + \sum_i y_i \log(f(k \cdot x_i)) + c \right) -$$

$$\lambda \left[-\Delta \sum_i x_i f'(k \cdot x_i) + \sum_i y_i x_i \frac{f'(k \cdot x_i)}{f(k \cdot x_i)} \right]$$

↳

• find critical points of \mathcal{L}

• find which critical point gives the min/max w.r.t k, λ

• that λ is the "Lagrange multiplier"

$$\text{STA} = \frac{1}{n} \sum_i y_i x_i \quad \therefore \hat{k} \text{ is a weighted STA w/ weights } \frac{f'}{f}$$

\hat{k} is same as ordinary STA if $f(z) = e^{az+b}$ i.e. is exponential

$$f'(z) = (a)e^{az+b}$$

$$\Rightarrow \frac{f'}{f} = a e^{az+b} / e^{az+b} = a$$

$$p(y_i | x_i, \theta) = \frac{1}{y_i!} [\Delta f(k \cdot x_i)]^{y_i} e^{-\Delta f(k \cdot x_i)}$$

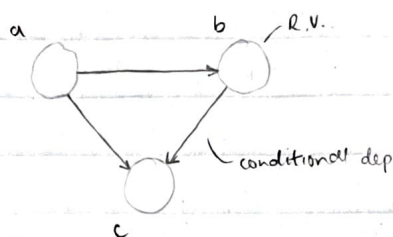
$$\text{if: } f(z) = e^{az+b}$$

$$p(y_i | x_i, \theta) = \frac{1}{y_i!} [\Delta e^{(a(k \cdot x_i) + b)}]^{y_i} e^{-\Delta e^{(a(k \cdot x_i) + b)}}$$

=

by the product rule (if prob):

$$p(a, b, c) = p(c|a, b) p(a, b) = p(c|a, b) p(b|a) p(a)$$



For a graph of K nodes, the joint dist.: $p(\mathbf{x}) = \prod_{k=1}^K p(x_k | pa_k)$
 where pa_k denotes the parents of x_k

example: polynomial regression

w : vector of polynomial coef.

t : vector of observed data $= (t_1, \dots, t_N)^T$

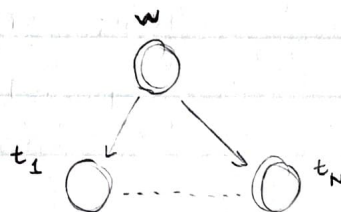
x : input data vector $= (x_1, \dots, x_N)^T$

σ^2 : variance

α : hyperparam. representing precision of the Gaussian prior over w

joint dist.:

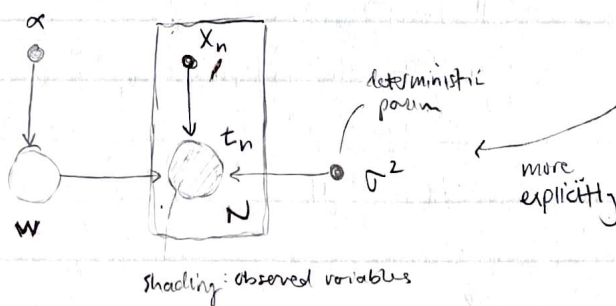
$$p(t, w) = p(w) \prod_{n=1}^N p(t_n | w)$$



re-write joint dist. more explicitly:

$$p(t, w | x, \alpha, \sigma^2) = p(w | \alpha) \prod_{n=1}^N p(t_n | w, x_n, \sigma^2)$$

\Downarrow equivalent to



posterior distribution of w :

$$p(w | T) \propto \underbrace{p(w)}_{\text{prior}} \prod_{n=1}^N \underbrace{p(t_n | w)}_{\text{likelihood}}$$

8.2 Conditional Independence

given: $p(a|b,c) = p(a|c)$ \therefore a is independent of b given c

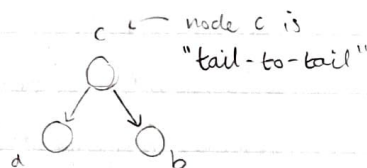
alternative expression: $p(a,b|c) = p(a|b,c)p(b|c)$
 $= p(a|c)p(b|c)$

" a & b are statistically independent given c "

\hookrightarrow shorthand: $a \perp\!\!\!\perp b | c$

example 1:

$$p(a,b,c) = p(a|c)p(b|c)p(c) \rightarrow$$



test if a & b are independent by marginalizing out c :

$$p(a,b) \stackrel{?}{=} p(a)p(b) = \sum_c p(a|c)p(b|c)p(c)$$

\rightarrow generally does not factorize $\Rightarrow a \not\perp\!\!\!\perp b | \emptyset$

"empty set"

instead show a & b are independent w.r.t c :

$$p(a,b|c) = p(a,b,c)/p(c) = p(a|c)p(b|c)$$

 $\Rightarrow a \perp\!\!\!\perp b | c$

$\therefore a$ & b are conditionally independent

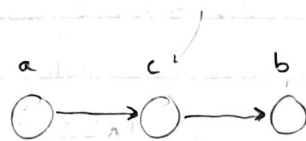
"conditioning on node c "blocks" the path b/w. a & b "

example 2:

given:

$$p(a, b, c) = p(a) p(c|a) p(b|c)$$

node c is "head-to-tail"



test if $a \perp b$ are independent by marginalizing over c :

$$\begin{aligned} p(a, b) &\stackrel{?}{=} p(a) p(b) = p(a) \sum_c p(c|a) p(b|c) \\ &= p(a) p(b|a) \end{aligned}$$

\rightarrow generally does not factorise to $p(a) p(b) \therefore a \not\perp b \mid \emptyset$

if condition on c :

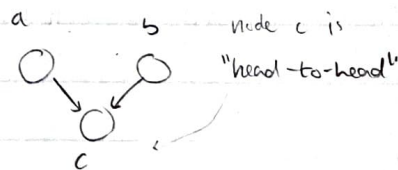
$$\begin{aligned} p(a, b|c) &= p(a, b, c) / p(c) \\ &= p(a) p(c|a) p(b|c) / p(c) \\ &= p(a|c) p(b|c) \end{aligned}$$

$\Rightarrow a \perp b \mid c$ " $a \perp b$ are conditionally independent w.r.t c "

example 3:

given:

$$p(a, b, c) = p(a) p(b) p(c|a, b) \rightarrow$$



test if $a \perp b$ are independent by marginalizing over c : $p(a, b) = p(a) p(b)$
 $\Rightarrow a \perp b \mid \emptyset$

test if $a \perp b$ are conditionally independent on c :

$$\begin{aligned} p(a, b|c) &= p(a, b, c) / p(c) \\ &= p(a) p(b) p(c|a, b) / p(c) \end{aligned}$$

\rightarrow does not generally factorize into $p(a) p(b) \therefore a \not\perp b \mid c$

- if c is not observed, then $a \perp b$ are independent
- if c is observed, then $a \perp b$ are dependent

9.2 Mixtures of Gaussians

Gaussian mixture: linear superposition of Gaussians:

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k N(\mathbf{x} | \mu_k, \Sigma_k)$$

 \mathbf{z} : K -dim binary r.v. where; $z_k \in \{0, 1\}$ and $\sum_k z_k = 1$

$$p(z_k = 1) = \pi_k \quad \text{where} \quad 0 \leq \pi_k \leq 1$$

$$\sum_k \pi_k = 1$$

$$\text{let } p(\mathbf{x}, \mathbf{z}) = p(\mathbf{z})p(\mathbf{x} | \mathbf{z})$$

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k} \quad \text{where } z_k \in \{0, 1\}$$

$$p(\mathbf{x} | z_k = 1) = N(\mathbf{x} | \mu_k, \Sigma_k)$$

$$p(\mathbf{x} | \mathbf{z}) = \prod_{k=1}^K N(\mathbf{x} | \mu_k, \Sigma_k)^{z_k} \quad \text{where } z_k \in \{0, 1\}$$

get $p(\mathbf{x})$ by marginalizing out \mathbf{z} from joint dist. $p(\mathbf{x}, \mathbf{z})$:

$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{z})p(\mathbf{x} | \mathbf{z})$$

sum over all values of \mathbf{z} :

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x} | \mathbf{z})$$

$$= \sum_{\mathbf{z}} \left[\prod_{k=1}^K \pi_k^{z_k} \cdot \prod_{k=1}^K N(\mathbf{x} | \mu_k, \Sigma_k)^{z_k} \right]$$

$$= \sum_{k=1}^K \pi_k N(\mathbf{x} | \mu_k, \Sigma_k)$$

* for each observation \mathbf{x}_n , there is a corresponding \mathbf{z}_n

9.2.2 EM for Gaussian mixtures

conditions that must be satisfied at a max. of the likelihood:

$$\text{log likelihood: } \ln p(\mathbf{X} | \pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left[\sum_{k=1}^K \pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k) \right]$$

$$\left(\frac{d}{d\mu} \ln p(\mathbf{X} | \pi, \mu, \Sigma) \right) \stackrel{!}{=} \text{set to } 0:$$

$$0 = - \sum_{n=1}^N \left[\frac{\pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_j \pi_j N(\mathbf{x}_n | \mu_j, \Sigma_j)} \Sigma_k (\mathbf{x}_n - \mu_k) \right]$$

$$\text{then } \mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

$$\text{where: } N_k = \sum_{n=1}^N \gamma(z_{nk}) \quad \text{"effective \# of data pts. assigned to class k"}$$

$$\gamma(z_{nk}) \equiv p(z_{nk} = 1 | \mathbf{x}_n) \quad \text{see eq. 9.13 for derivation}$$

\therefore mean of k^{th} Gaussian μ_k is the mean of all the data points weighted by their probability of being in class k

$$\frac{d}{d\Sigma} \ln p(\mathbf{X} | \pi, \mu, \Sigma) \stackrel{!}{=} \text{set to } 0:$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T$$

\therefore covariance matrix for Gaussian k (Σ_k) is a weighted average of covariance of all data points

$$\frac{d}{d\pi} \ln p(\mathbf{X} | \pi, \mu, \Sigma) \stackrel{!}{=} \text{set to } 0: \quad \pi_k = \frac{N_k}{N}$$

posterior prob. given \mathbf{x}
for class k

\therefore the mixing coef. for class k is the average "responsibility" which the component takes for explaining the data points

9.3. An alternative view of EM

* goal of EM: find the ^{max. likelihood} ML solutions for models having latent variables

X : all observed data; n^{th} row represents x_n^T

Z : all latent vars.; n^{th} row represents z_n^T

θ : all model parameters

log likelihood:

$$\ln p(X|\theta) = \ln \left[\sum_Z p(X, Z|\theta) \right]$$

complete data set: $\{X, Z\}$ as if we knew the latent variable values

incomplete data set: X alone

Expectation step: E

1. use the current values θ^{old} to find the posterior dist. of the latent variables by $p(Z|X, \theta^{\text{old}})$
2. use this posterior dist. to find the expectation of the complete-data

log likelihood:

$$Q(\theta, \theta^{\text{old}}) = \sum_Z p(Z|X, \theta^{\text{old}}) \ln p(X, Z|\theta)$$

Maximization step: M

3. determine the revised parameter estimate θ^{new} by maximizing $Q(\theta, \theta^{\text{old}})$

$$\theta^{\text{new}} = \operatorname{argmax}_{\theta} Q(\theta, \theta^{\text{old}})$$

1. Generalized linear models recap

2. Conditional independence & graphical models

· when talking about d-sep., "observing" a value for a node = conditioning on the node

3. Mixture models and the EM algorithm

· latent var. formalization of Gaussian mix. model is done so can use E.M.

Exercise #2 notes

2/20/2022

[black \rightarrow left white \rightarrow right]

DATA

"activity" : [timepoint \times neuron] spikecounts

"cue" : experimental cue ; 0 = black , 1 = white

"cho" : choice ; 0 = left , 1 = right

"corr" : correct choice ; 0 = incorrect , 1 = correct

"prev_corr" : previous outcome ; " " " " " "

Section 2.3 : Fitting Poisson GLM

$$\log(y) = w_0 + w_1 x_1 + \dots \quad w: \text{weights} \quad x: \text{data/features}$$