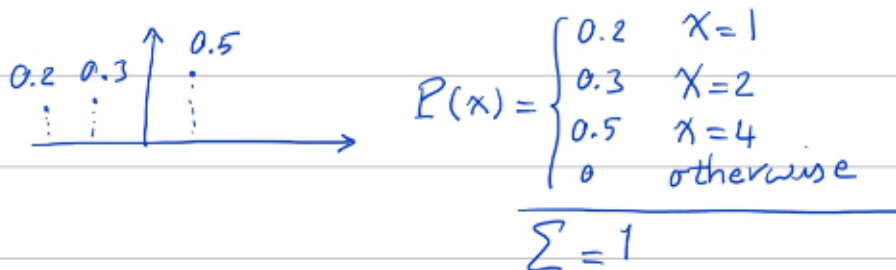
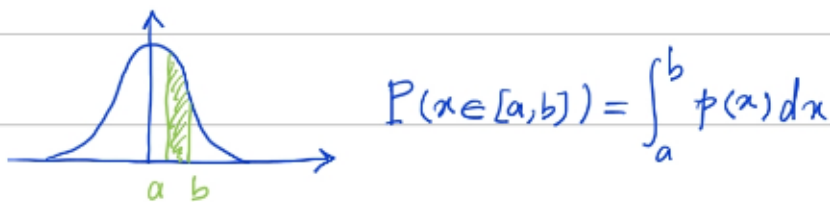


Discrete vs Continuous Distributions



- In discrete approach $P(x)$ is directly given



- In continuous approach, the most convenient way is by defining Probability Distribution Function (PDF)

Dependencies

- The two ran variables are considered independent if:

$$P(X, Y) = P(X) P(Y)$$

↘ joint
↑ Marginals

Conditional Probability

$$P(X|Y) = \frac{P(X, Y)}{P(Y)}$$

↘ Conditional
↘ Marginal

Dependencies

• The two random variables are considered independent if:

$$P(X, Y) = P(X) P(Y)$$

↘ joint ↗ Marginals

Conditional Probability

$$P(X|Y) = \frac{P(X, Y)}{P(Y)}$$

↙ Conditional ↘ Marginal

Trick 1: Chain Rule

$$P(X, Y) = P(X|Y) P(Y)$$

$$P(X, Y, Z) = P(X|Y, Z) P(Y|Z) P(Z)$$

$$P(X_1, \dots, X_N) = \prod_{i=1}^N P(X_i | X_1, \dots, X_{i-1})$$

Note: The intersection (joint probability) is associative/

communicative: Therefore:

$$P(X, Y, Z) = P(X|Y, Z) P(Y|Z) P(Z)$$

$$= P(X|Y, Z) P(Z|Y) P(Y)$$

$$= P(Y|X, Z) P(Z|X) P(X)$$

⋮

Trick 2: Sum Rule

Note: The intersection (joint probability) is associative

commutative: Therefore:

$$P(X, Y, Z) = P(X|Y, Z) P(Y|Z) P(Z)$$

$$= P(X|Y, Z) P(Z|Y) P(Y)$$

$$= P(Y|X, Z) P(Z|X) P(X)$$

⋮

Trick 2: Sum Rule

$$p(x) = \int_{-\infty}^{\infty} p(x, y) dy$$

Annotations:

- $p(x)$: marginal distribution
- $p(x, y)$: joint probability
- y : random variable

Bayes Theorem

$$P(\theta|X) = \frac{P(\theta, X)}{P(X)} = \frac{P(X|\theta) P(\theta)}{P(X)}$$

Annotations:

- $P(\theta|X)$: posterior
- $P(\theta, X)$: joint probability
- $P(X|\theta)$: likelihood
- $P(\theta)$: prior
- $P(X)$: evidence

Frequentist

vs

Bayesian

objective

subjective

considers randomness of events

tries to predict behavior based on evidences

$\begin{cases} \theta \text{ is fixed} \\ X \text{ is random} \end{cases}$

$\begin{cases} \theta \text{ is random} \\ X \text{ is fixed} \end{cases}$

considers randomness of events

$\begin{cases} \theta \text{ is fixed} \\ X \text{ is random} \end{cases}$

uses no previous behavior based on evidences

$\begin{cases} \theta \text{ is random} \\ X \text{ is fixed} \end{cases}$

works only when the number of data points is much bigger than the number of parameters

$$|X| \gg |\theta|$$

works for arbitrary number of data points

for any $|X|$

For training, it uses Maximum Likelihood

$$\hat{\theta} = \arg \max_{\theta} P(X|\theta)$$

for training, it uses Bayes Theorem

$$P(\theta|X) = \frac{P(X|\theta)P(\theta)}{P(X)}$$

Training

$$P(\theta|X_{tr}, y_{tr}) = \frac{P(y_{tr}|X_{tr}, \theta)P(\theta)}{P(y_{tr}|X_{tr})}$$

Prediction

$$P(y_{ts}|X_{ts}, X_{tr}, y_{tr}) = \int P(y_{ts}|X_{ts}, \theta)P(\theta|X_{tr}, y_{tr})d\theta$$

Regularization

Regularizer

$$P(\theta|X) = \frac{P(X|\theta)P(\theta)}{P(X)}$$

online learning

Likelihood

Prior

$$P_k(\theta) = P(\theta|x_k) = \frac{P(x_k|\theta)P_{k-1}(\theta)}{P(x_k)}$$

new prior

posterior

Draft Page

-proof of training formula below:

$$P(y_{tr}|X_{tr}, \theta)P(\theta)$$

- proof of ^u training formula below:

$$* P(\theta | X_{tr}, y_{tr}) = \frac{P(y_{tr} | X_{tr}, \theta) P(\theta)}{P(y_{tr} | X_{tr})}$$

$$P(\theta | X_{tr}, y_{tr}) \underset{\text{Bayes Rule}}{=} \frac{P(\theta, X_{tr}, y_{tr})}{P(X_{tr}, y_{tr})}$$

$$\xrightarrow{\text{Chain Rule}} = \frac{P(y_{tr} | \theta, X_{tr}) P(X_{tr} | \theta) P(\theta)}{P(y_{tr} | X_{tr}) P(X_{tr})} = \frac{P(y_{tr} | X_{tr}, \theta) P(\theta)}{P(y_{tr}, X_{tr})}$$

Proof of equation of prediction

$$P(y_{ts} | X_{ts}, X_{tr}, y_{tr}) \underset{\text{Sum Rule}}{=} \int P(y_{ts} | X_{ts}, X_{tr}, y_{tr}, \theta) d\theta$$

$$P(y_{ts} | X_{ts}, X_{tr}, y_{tr}, \theta) \underset{\text{Bayes Rule}}{=} \frac{P(y_{ts}, X_{ts}, X_{tr}, y_{tr}, \theta)}{P(X_{ts}, X_{tr}, y_{tr}, \theta)}$$
$$\rightarrow =$$

How to define a model?

Bayesian network

Nodes: random variables

Edges: direct impact

Model: joint probability over all variables

$$P(X_1, \dots, X_n) = \prod_{k=1}^n P(X_k | \text{Pa}(X_k))$$

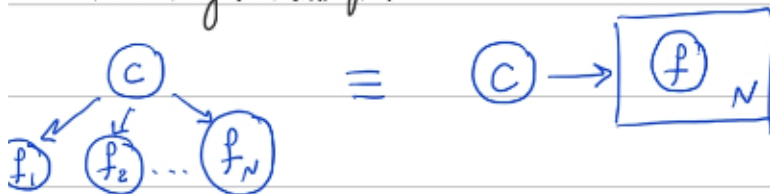
→ parents



$$P(S, R, G) = P(G|S, R)P(S|R)P(R)$$

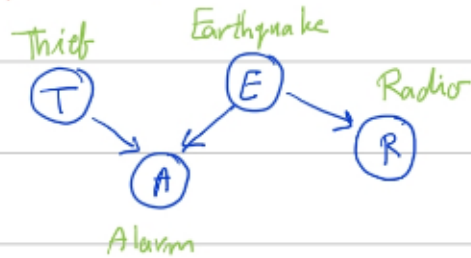
Naïve Bayes classifier

→ Plate Notation



$$P(C, f_1, f_2, \dots, f_N) = P(C) \prod_{i=1}^N P(f_i | C)$$

Example: Thief & Alarm



$$P(t,a,e,r) = ? = P(t) P(e) P(a|t,e) P(r|e)$$

Priors	
$P(T=1)$	10^{-3}
$P(E=1)$	10^{-2}

$P(A=1 T,E)$	$E=0$	$E=1$
$T=0$	0	1/10
$T=1$	1	1

$P(R E)$	
$E=0$	0
$E=1$	1/2

* a notation simplification: $T=1 \rightarrow T$
 $T=0 \rightarrow \bar{T}$

$$P(T|A) = \frac{P(T,A)}{P(A)} = \frac{P(T,A,E) + P(T,A,\bar{E})}{P(T,A,E) + P(T,A,\bar{E}) + P(\bar{T},A,E) + P(\bar{T},A,\bar{E})} \quad \rightarrow \text{Sum Rule}$$

$$\bullet P(T,A,E) = P(A|T,E) P(T) P(E) = 10^{-3} \times 10^{-2} \times 1 = 10^{-5}$$

$$\bullet P(\bar{T},A,\bar{E}) = P(A|\bar{T},\bar{E}) P(\bar{T}) P(\bar{E}) = 0$$

$$\bullet P(T,A,\bar{E}) = P(A|T,\bar{E}) P(T) P(\bar{E}) = 1 \times 10^{-3} \times 0.99 = 9.9 \times 10^{-4}$$

$$\therefore P(T|A) \simeq 50\%$$

$P(T|A,R) = ?$ what is the probability of a Thief in a house if you hear an alarm & radio report

$$P(T|A,R) = \frac{P(T,A,R)}{P(A,R)} = \frac{P(T,A,R,E) + P(T,A,R,\bar{E})}{P(A,R,T,E) + P(A,R,T,\bar{E}) + P(A,R,\bar{T},E) + P(A,R,\bar{T},\bar{E})}$$

* The earthquake parameter 'E' is added because the event depends on 'E' as a parent. *

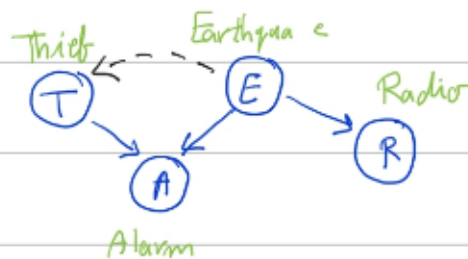
$$\circ P(T,A,R,E) = P(A|T,E) P(R|E) P(E)$$

$$\circ P(T,A,R,\bar{E}) = P(A|T,\bar{E}) P(R|\bar{E}) P(\bar{E}) = 1 \times 0 \times 0.99 = 0$$

$$\therefore \Rightarrow P(T|A,R) \simeq 1\%$$

* If the results do not match with expectations, the model should be modified. For example:

- There are more Thieves in an event of an earthquake



LINEAR REGRESSION

Univariate Normal

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Multivariate Normal

$$\mathcal{N}(x|\mu, \Sigma) = \frac{1}{\sqrt{|2\pi\Sigma|}} \exp\left[-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right]$$

$$\mathbb{E}X = \mu \quad \text{Cov}[X] = \Sigma$$

mean vector

Covariant Matrix

Σ matrix representations

$$\Sigma = \begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix}$$



Full
* params: $\frac{D(D+1)}{2}$

$$\Sigma = \begin{pmatrix} 3 & 0 \\ 0 & 1 \end{pmatrix}$$



Diagonal
* params: D

$$\Sigma = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$$



Spherical
* params: 1

approximation

approximation

Linear Regression

$$L(\omega) = \sum_{i=1}^N (\omega^T x_i - y_i)^2 = \|\omega^T X - y\|^2 \rightarrow \min_{\omega}$$

$$\hat{\omega} = \arg \min_{\omega} L(\omega)$$

Bayesian Model



$$P(\omega, y|x) = P(y|x, \omega) P(\omega)$$

$$P(y|x, \omega) = \mathcal{N}(y|\omega^T x, \sigma^2 I) : \text{assuming normal distribution}$$

$$P(\omega, y|X) = P(y|X, \omega) P(\omega)$$

$$P(y|X, \omega) = \mathcal{N}(y|\omega^T X, \sigma^2 I) : \text{assuming normal distribution}$$

$$P(\omega) = \mathcal{N}(\omega|0, \gamma^2 I)$$

- Calculate the posterior probability

$$P(\omega|y, X) = \frac{P(y, \omega|X)}{P(y|X)} \rightarrow \max_{\omega}$$

depends on ω
independent of ω

$$\rightarrow P(y, \omega|X) = P(y|X, \omega) P(\omega) \rightarrow \max_{\omega} \rightarrow \text{take log}$$

$$\log P(y|X, \omega) + \log P(\omega) =$$

$$\log C_1 \exp\left(-\frac{1}{2} (y - \omega^T X)^T [\sigma^2 I]^{-1} (y - \omega^T X)\right) +$$

$$\log C_2 \exp\left(-\frac{1}{2} \omega^T [\gamma^2 I]^{-1} \omega\right)$$

$$= \frac{1}{2\sigma^2} \underbrace{(y - \omega^T X)^T (y - \omega^T X)}_{\|y - \omega^T X\|^2} - \frac{1}{2\gamma^2} \underbrace{\omega^T \omega}_{\|\omega\|^2} \rightarrow \max_{\omega} \quad |x - 1 \times 2\sigma^2$$

$$\Rightarrow \underbrace{\|y - \omega^T X\|^2}_{\text{sum of squares}} + \underbrace{\lambda \|\omega\|^2}_{L_2 \text{ Regularizer}} \rightarrow \min_{\omega}$$

Exam Tips:

2-1: When b and c are independent:

$$p(a|b) = \int p(a, c|b) dc$$

$$p(a, c|b) = \frac{p(a, b, c)}{p(b)} = \frac{p(a|b, c) p(b|c) p(c)}{p(b)}$$

if b and c are independent: $p(b|c) = p(b)$

$$\Rightarrow p(a|b) = \int p(a|b, c) p(c) dc$$

2-2

$$\textcircled{1} \quad p(a, c|b) = \frac{p(a, b, c)}{p(b)} \quad , \quad \textcircled{2} \quad p(c|a, b) = \frac{p(a, b, c)}{p(a, b)}$$

$$\Rightarrow \frac{\textcircled{1}}{\textcircled{2}} = \frac{p(a, b)}{p(b)} = p(a|b)$$

Analytical Inference

$$\text{Posterior distribution: } P(\theta|X) = \frac{P(X|\theta)P(\theta)}{P(X)}$$

Likelihood \swarrow prior \swarrow
evidence \swarrow

• What is evidence $P(X)$?

* Maximum a posteriori principle:

- We try to find the value of parameters that maximizes the posterior probability:

$$\theta_{\text{MP}} = \arg \max_{\theta} P(\theta|X) = \arg \max_{\theta} \frac{P(X|\theta)P(\theta)}{P(X)}$$

$$\theta_{\text{MP}} = \arg \max_{\theta} P(X|\theta)P(\theta)$$

\hookrightarrow avoid computing the evidence

* problems

- Not invariant to reparameterization

- Can't be used as a prior

- MAP is a solution to $L(\theta) = \mathbb{I}(\theta \neq \theta^*) \rightarrow \min_{\theta}$

- can't compute credible regions

Conjugate Distributions (another method to avoid computing evidence)

Fixed by model

our own choice

$$P(\theta|x) = \frac{P(x|\theta) P(\theta)}{P(x)}$$

Fixed by data

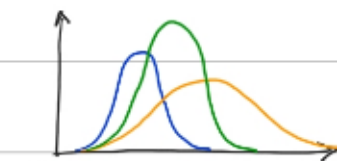
$P(\theta)$ is conjugate to $P(x|\theta)$:

$\mathcal{A}(v)$

$$\mathcal{A}(v') \rightarrow P(\theta|x) = \frac{P(x|\theta) P(\theta)}{P(x)}$$

Example: Two Gaussians

$\mathcal{N}(x|\theta, \sigma^2)$



$\mathcal{N}(\theta|m, s^2)$

$$P(\theta|x) = \frac{P(x|\theta) P(\theta)}{P(x)}$$

$\mathcal{N}(\theta|a, b^2)$

Example: Gamma distribution

$$\Gamma(y|a,b) = \frac{b^a}{\Gamma(a)} y^{a-1} e^{-by}$$

$$y, a, b > 0 \quad \uparrow \quad \Gamma(n) = (n-1)!$$

mean: $E(y) = a/b$

$$\text{Mode}[y] = \frac{a-1}{b}$$

$$\text{Var}[y] = a/b^2$$

* Gamma distribution is conjugate to Normal distribution

with respect to precision (inverse of variance $\gamma = \frac{1}{\sigma^2}$)

$$\mathcal{N}(x|\mu, \gamma^{-1}) = \frac{\sqrt{\gamma}}{\sqrt{2\pi}} e^{-\frac{\gamma(x-\mu)^2}{2}} \propto \gamma^{1/2} e^{-b\gamma}$$

assume: $p(\gamma) \propto \gamma^{a-1} e^{-b\gamma} = \Gamma(\gamma|a,b)$

let's check: $p(\gamma|x) \propto p(x|\gamma)p(\gamma)$

$$p(\gamma|x) \propto \left(\gamma^{1/2} e^{-\frac{\gamma(x-\mu)^2}{2}} \right) \cdot \left(\gamma^{a-1} e^{-b\gamma} \right)$$

$$p(\gamma|x) \propto \gamma^{1/2+a-1} e^{-\gamma(b + \frac{(x-\mu)^2}{2})}$$

$$\Rightarrow p(\gamma, x) = \Gamma(a + 1/2, b + \frac{(x-\mu)^2}{2})$$

Example: Beta distribution (good for modeling final support)

$$B(x|a, b) = \frac{1}{B(a, b)} x^{a-1} (1-x)^{b-1}$$

$$x \in [0, 1], a, b > 0 \quad \uparrow \quad \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}$$

$$\text{mean: } E x = \frac{a}{a+b}$$

$$\text{Mode}[x] = \frac{a-1}{a+b-2}$$

$$\text{Var}[x] = \frac{ab}{(a+b)^2(a+b-1)}$$

* Beta distribution is conjugate to Bernoulli likelihood

— Bernoulli likelihood

$$p(x|\theta) = \theta^{N_1} (1-\theta)^{N_0}$$

$$\text{and: } p(\theta) = B(\theta|a, b) \propto \theta^{a-1} (1-\theta)^{b-1}$$

$$p(\theta|x) \propto p(x|\theta) p(\theta)$$

$$p(\theta|x) \propto \theta^{N_1} (1-\theta)^{N_0} \cdot \theta^{a-1} (1-\theta)^{b-1}$$

$$p(\theta|x) \propto \theta^{N_1+a-1} (1-\theta)^{N_0+b-1} : \text{a Beta dist. func.}$$

$$\Rightarrow p(\theta|x) = B(N_1+a, N_0+b)$$

Summary

Pros:

- Exact posteriors
- Easy for online learning

Cons:

- Conjugate prior may be inadequate

Latent Variable Models

Pros: • Simpler models (less edges)

• fewer parameters

• Latent variables are sometimes meaningful

Cons: • Harder to work with

— Why using Bayesian method instead of standard regression approaches sometimes?

• There are some missing values

• To quantify uncertainty in predictions

— Latent Variable Models

• The connections of Bayesian model could be too complex.

• Therefore the table of probabilities could be very large

• The "Intelligence" is introduced to model/capture some features of each component.

* An Example:

$$p(x_1, \dots, x_5) = \sum_{I=1}^{100} p(x_1, \dots, x_5 | I) p(I)$$

↑ Sum over the number of assigned intelligence

$$= \sum_{I=1}^{100} p(x_1 | I) \dots p(x_5 | I) p(I)$$

Five tables only are needed

Probabilistic Clustering

- Soft Clustering

$p(\text{cluster id} | x)$ instead of $\text{cluster id} = f(x)$

- It assumes that each point belongs to all clusters, but with different probabilities.

- Benefits of using...

- Handling missing values naturally
- Hyperparameter tuning
- Generating new data points

* Some Examples:

$$\rightarrow p(x | \theta) = \mathcal{N}(x | \mu, \Sigma)$$

$$\theta = \{\mu, \Sigma\}$$

- In case one Gaussian is not able to cover all data points one alternative is to use multiple distributions:

Gaussian Mixture Model (GMM)

$$p(x | \theta) = \pi_1 \mathcal{N}(x | \mu_1, \Sigma_1) + \pi_2 \mathcal{N}(x | \mu_2, \Sigma_2) + \pi_3 \mathcal{N}(x | \mu_3, \Sigma_3)$$

$$\theta = \{\pi_1, \pi_2, \pi_3, \mu_1, \mu_2, \mu_3, \Sigma_1, \Sigma_2, \Sigma_3\}$$

How to train a GMM (estimation of params) ?

- By finding maximum likelihood: $\max_{\theta} p(X|\theta)$

$$\max_{\theta} p(X|\theta) = \prod_{i=1}^N p(x_i|\theta)$$

$$= \prod_{i=1}^N (\pi_1 \mathcal{N}(x_i|\mu_1, \Sigma_1) + \pi_2 \mathcal{N}(x_i|\mu_2, \Sigma_2) + \dots)$$

N : number of data points

&

$$\pi_1 + \pi_2 + \pi_3 = 1; \pi_k \geq 0; k=1, 2, 3$$

* For this optimization problem, tools like tensorflow can be used.

* The covariance matrix Σ is not arbitrary:

It should be positive semi definite : $\Sigma \succcurlyeq 0$

* Why not to use SGD in this case?

* It is hard to follow some constraints

* Expectation Maximization Algorithm, which can exploit the structure of the program, sometimes is much more faster and efficient.

Training GMM (Gaussian Mixture Model)

Q: How to do better than SGD?

$$p(x|\theta) = \pi_1 \mathcal{N}(x|\mu_1, \Sigma_1) + \pi_2 \mathcal{N}(x|\mu_2, \Sigma_2) + \pi_3 \mathcal{N}(x|\mu_3, \Sigma_3) \quad (\text{eq. 1})$$

* Assign a latent variable "t" for each data point x.



* For this problem, $t=1, 2, 3$ (three values), showing that which Gaussian this particular data point came from (which is unknown)

* it is reasonable to assume that "t" has prior distribution π .

* if we know that a data point comes from Gaussian number C:

$$p(t=C|\theta) = \pi_C \leftarrow \text{prior}$$

* The density of data point, given the cluster ~~x~~:

$$p(x|t=C, \theta) = \mathcal{N}(x|\mu_C, \Sigma_C)$$

* Rule of sum for probabilities:

$$p(x|\theta) = \sum_{C=1}^3 p(x|t=C, \theta) p(t=C|\theta) \quad (\text{eq. 2})$$

↖ Marginalizing "t"

* eq. 1 is exactly the same as eq. 2.

Expectation Maximization

* How to train this latent variable model?

— How to estimate θ ?

— From training data set, if "t" is known (for example):

$$p(x|t=1, \theta) = \mathcal{N}(x|\mu_1, \sigma_1^2)$$

$$\mu_1 = \frac{\sum_{\text{blue}} x_i}{\# \text{ of blue points}}, \quad \sigma_1^2 = \frac{\sum_{\text{blue}} (x_i - \mu_1)^2}{\# \text{ of blue points}}$$

$$\mu_1 = \frac{\sum_i p(t_i=1|x_i, \theta) x_i}{\sum_i p(t_i=1|x_i, \theta)}$$

$$\sigma_1^2 = \frac{\sum_i p(t_i=1|x_i, \theta) (x_i - \mu_1)^2}{\sum_i p(t_i=1|x_i, \theta)}$$

← can also be
calculated from
posterior

* How can we choose the best run among several

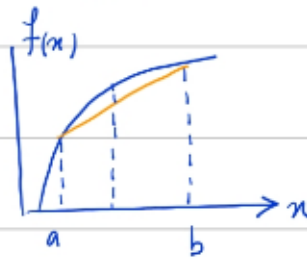
training attempts with different random initialization?

— choose the one with the highest training log-likelihood

— choose the one with the highest validation log-likelihood...

General Form of Expectation Maximization

* Concave Functions



for any a, b, α : $f(\alpha a + (1-\alpha)b) \geq \alpha f(a) + (1-\alpha)f(b)$

$$0 \leq \alpha \leq 1$$

* Jensen's inequality

If " f " is a concave function and

$$\alpha_1 + \alpha_2 + \alpha_3 = 1 ; \alpha_k \geq 0$$

$$\Rightarrow f(\underbrace{\alpha_1 a_1 + \alpha_2 a_2 + \alpha_3 a_3}_{E_{p(t)} t}) \geq \underbrace{\alpha_1 f(a_1) + \alpha_2 f(a_2) + \alpha_3 f(a_3)}_{E_{p(t)} f(t)}$$

$$E_{p(t)} t$$

$$E_{p(t)} f(t)$$

$$p(t=a_1) = \alpha_1$$

$$p(t=a_2) = \alpha_2$$

$$p(t=a_3) = \alpha_3$$

$$\Rightarrow f(E_{p(t)} t) \geq E_{p(t)} f(t)$$

Kullback-Leibler Divergence

↳ It's a way to measure difference between two probabilistic distributions.

$$KL(q||p) = \int q(x) \log \frac{q(x)}{p(x)} dx$$

\int is Expected value of $q(x)$ a measure of diff

$$* KL(q||p) \neq KL(p||q)$$

$$* KL(q||q) = 0$$

$$* KL(q||p) \geq 0$$

proof: $-KL(q||p) = E_q(-\log \frac{q}{p}) = E(\log \frac{p}{q})$

$$\leq \log(E_q \frac{p}{q}) = \log \int q(x) \frac{p(x)}{q(x)} dx = 0$$

General form of Expectation Maximization

* to find out $\max_{\theta} p(X|\theta)$, max of likelihood

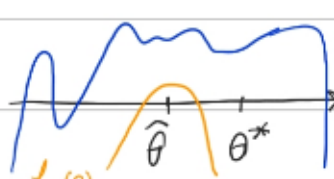
\nwarrow # of data points

$$\max_{\theta} \log p(X|\theta) = \log \prod_{i=1}^N p(x_i|\theta)$$

$$\log p(X|\theta) = \sum_{i=1}^N \log p(x_i|\theta)$$

$$(eq.1) = \sum_{i=1}^N \log \sum_{c=1}^3 p(x_i, t_i=c|\theta) \geq L(\theta)$$

a lower bound - Jensen's inequality \nearrow

e.g.  * one lower bound is not enough, so we define a family of lower bounds.

$$(eq.1) = \sum_{i=1}^N \log \sum_{c=1}^3 \frac{q(t_i=c)}{q(t_i=c)} p(x_i, t_i=c|\theta)$$

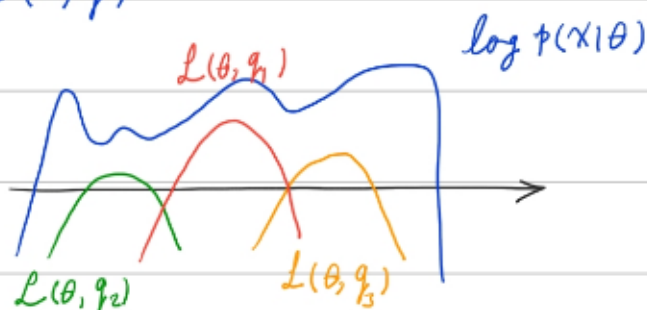
$$\boxed{\log \left(\sum_c \alpha_c v_c \right) \geq \sum_c \alpha_c \log(v_c)} \quad \text{Jensen's inequality}$$

$$\alpha_c: q(t_i=c)$$

$$v_c: \frac{p(x_i, t_i=c|\theta)}{q(t_i=c)}$$

$$\Rightarrow (eq.1) \geq \sum_{i=1}^N \sum_{c=1}^3 q(t_i=c) \log \frac{p(x_i, t_i=c|\theta)}{q(t_i=c)}$$

$$= L(\theta, q)$$



* Iteration Steps:

— For θ_k , find $L(\theta, q)$ which is maximum at this current point θ_k / find q such that the value of a lower bound of the point θ_k and q is maximum.
(check PDF of slides)

E-Step $q^{k+1} = \arg \max_q L(\theta^k, q)$

M-Step $\theta^{k+1} = \arg \max_{\theta} L(\theta, q^{k+1})$

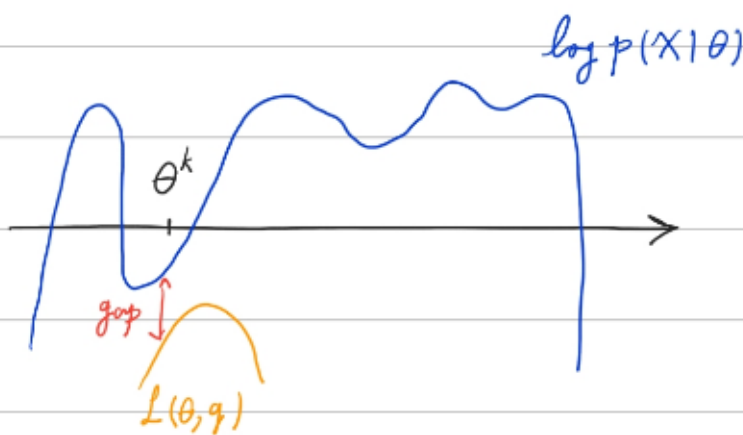
E-Step Details

$$\log p(X|\theta) \geq \mathcal{L}(\theta, q)$$

q : itself is a distribution

E-Step: $\max_q \mathcal{L}(\theta^k, q)$

1) choose the one which has the highest value at θ_k .



gap:

$$\log p(X|\theta) - \mathcal{L}(\theta, q) = \sum_i K \mathcal{L}(q(t_i) || p(t_i | x_i, \theta))$$

proof:

- we assume that the data-set consists of the objects that are independent of given parameters.

$$= \sum_{i=1}^N \log p(x_i | \theta) - \sum_{i=1}^N \log \sum_{c=1}^3 \underbrace{q(t_i=c)}_{\text{orange}} \frac{p(x_i, t_i=c | \theta)}{q(t_i=c)}$$

$$= \sum_{i=1}^N \left(\log p(x_i | \theta) \cdot \underbrace{\sum_c q(t_i=c)}_1 - \sum_{c=1}^3 q(t_i=c) \cdot \log \dots \right)$$

$$= \sum_{i=1}^N \sum_c q(t_i=c) \left(\log p(x_i|\theta) - \log \frac{p(x_i, t_i=c|\theta)}{q(t_i=c)} \right)$$

$$= \sum_{i=1}^N \sum_c q(t_i=c) \log \frac{p(x_i|\theta) q(t_i=c)}{p(x_i, t_i=c|\theta)}$$

\downarrow
 $p(t_i=c|x_i, \theta) p(x_i|\theta)$

$$= \sum_{i=1}^N \sum_c q(t_i=c) \log \frac{q(t_i=c)}{p(t_i=c|x_i, \theta)}$$

$\underbrace{\hspace{10em}}_{\text{KL}(q(t_i) || p(t_i|x_i, \theta))}$

$$= \sum_{i=1}^N \text{KL}(q(t_i) || p(t_i|x_i, \theta)) \quad (\text{Eq. 2})$$

$$\text{GAP} = \log p(X|\theta) - \underbrace{\underbrace{\mathcal{L}(\theta, q)}_{\text{max } q}}_{\text{min } q}$$

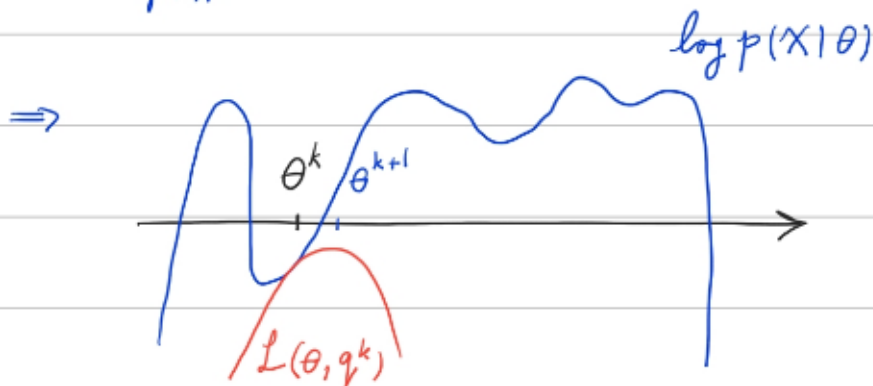
$\underbrace{\hspace{10em}}_{\text{min } q}$

$$\Rightarrow \min_q (\text{Eq. 2}) \geq 0 \Rightarrow q(t_i) = p(t_i|x_i, \theta)$$

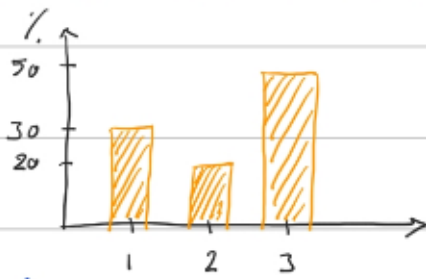
For global optimal

* E-Step

$$\arg \max_{q(t_i)} \mathcal{L}(\theta^k, q) = p(t_i|x_i, \theta)$$



Example: EM for Discrete Mixture, E-Step



* Assumption:

- probability of x_i : mixture of two distributions

$$p(x_i) = \gamma p_1(x_i) + (1-\gamma) p_2(x_i)$$

- let's say that:

	1	2	3
p_1	α	$1-\alpha$	0
p_2	0	$1-\beta$	β

- objective: to estimate α, β , and γ via EM method.

1) parameters initialization:

$$\alpha_0 = \beta_0 = \gamma_0 = 0.5$$

2) Define the latent variable t_i for each x_i



- t_i can just take two values

$$\text{e.g.: } p(t_i=1) = \gamma, \quad p(t_i=2) = 1-\gamma, \dots$$

$$\Rightarrow p(x_i | t_i=2) = p_2(x_i)$$

- E-Step:

Finding out the posterior distribution on the latent variable t_i :

$$\Rightarrow q(t_i=c) = p(t_i=c | x_i)$$

Start with:

$$\begin{aligned} p(t_i=1 | x_i=1) &= \frac{p(x_i=1 | t_i=1) p(t_i=1)}{p(x_i=1 | t_i=1) p(t_i=1) + p(x_i=1 | t_i=2) p(t_i=2)} \\ &= \frac{\gamma}{\gamma + (1-\gamma)} = 1 \end{aligned}$$

$$\begin{aligned} p(t_i=1 | x_i=3) &= \frac{p(x_i=3 | t_i=1) p(t_i=1)}{p(x_i=3 | t_i=1) p(t_i=1) + p(x_i=3 | t_i=2) p(t_i=2)} \\ &= 0 \end{aligned}$$

$$p(t_i=1 | x_i=2) = \dots = \frac{0.5 \times 0.5}{0.5 \times 0.5 + 0.5 \times 0.5} = 0.5$$

Example: EM for discrete mixture: M-Step

Objective:

$$\max_{\alpha, \beta, \gamma} \sum_{i=1}^N \mathbb{E}_{q(t_i)} \log \overbrace{p(x_i | t_i) p(t_i)}^{p(x_i, t_i)}$$

Summary of E-Step:

$$q(t_i=1) = p(t_i=1 | x_i) = \begin{cases} 1 & ; x_i=1 \\ 0.5 & ; x_i=2 \\ 0 & ; x_i=3 \end{cases}$$

$$q(t_i=2) = 1 - q(t_i=1)$$

$$= \sum_{i=1}^N q(t_i=1) \log p(x_i | t_i=1) \gamma + \sum_{i=1}^N q(t_i=2) \log p_2(x_i) (1-\gamma)$$

* assumption: $N_1=30$, $N_2=20$, $N_3=60$

$$\begin{aligned} &= 30 \cdot \underbrace{p(t_i=1 | x_i=1)}_1 \log \alpha \gamma + 20 \cdot 0.5 \cdot \log(1-\alpha) \gamma \\ &\quad + \cancel{60 \cdot 0 \cdot \log 0} + 30 \cdot \cancel{p(t_i=2 | x_i=1)}_0 \dots + 20 \cdot 0.5 \cdot \log(1-\beta)(1-\gamma) \\ &\quad + 60 \cdot 1 \cdot \log \beta (1-\gamma) \end{aligned}$$

maximizing wrt α : \rightarrow just consider terms depending on α

$$f_1(\alpha) = 30 \log \alpha + 10 \log(1-\alpha) + \text{const}(\alpha) \Rightarrow \max f_1(\alpha)$$

$$\frac{\partial f_1(\alpha)}{\partial \alpha} = 0 \Rightarrow 30 \frac{1}{\alpha} + 10 \frac{(-1)}{1-\alpha} = 0 \Rightarrow \dots$$

$$\text{summary: } \alpha = \frac{3}{4}, \beta = \frac{6}{7}, \gamma = \frac{4}{11}$$

meaning:

	1	2	3
P_1	α	$1-\alpha$	0
P_2	0	$1-\beta$	β

General EM for GMM

Objective: How to apply the EM algorithm to some concrete latent variable models.