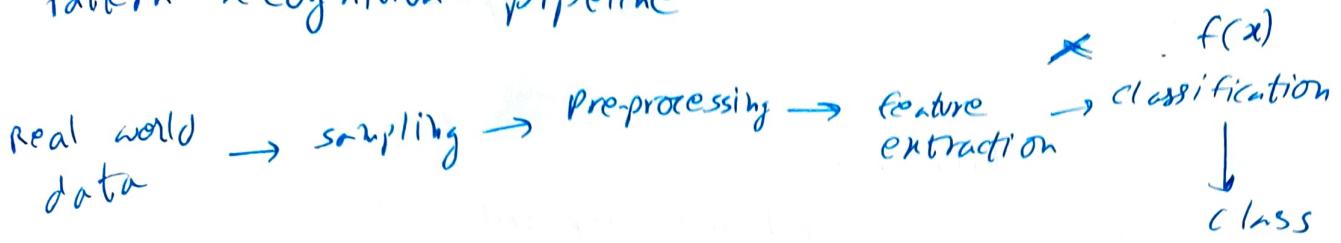


pattern recognition pipeline



$X \rightarrow$ feature vector

this is classified using a classifier or function $f(x)$ which op's the class y after classification.

Assumption: similar objects map to similar features.

supervised learning

Training sample is a tuple (x_i, y_i) .

\downarrow
feature vector \hookrightarrow expected class/desired prediction.

Unsupervised learning

work with unlabeled data. $(x_i) \rightarrow$ operates only on i/p x_i .

- summary of distribution

parameters

- choice of param. is training

hyper-parameters

- a parameter that parametrizes the choice of parameters.

- d is hyper-parameter

linear regression

$$y = \beta^T \tilde{x} = \sum_{i=0}^d \beta_i \tilde{x}_i$$

$\tilde{x}_i = (1, x_1, \dots, x_{d-1})^T \rightarrow$ plane defined in some space.
 \hookrightarrow offset

Flexibility requires more localization, more parameters A are less robust.

Model selection prob: trade-off w/ choice of parameters/hyperparam.

and flexibility.
Local models perform poorly in high dimensional spaces.

Prob. distributions are used to represent data in feature spaces.

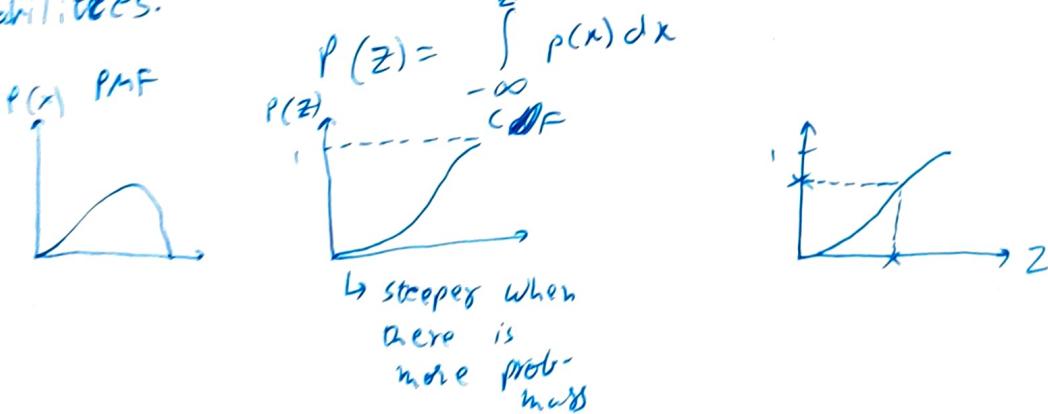
operations on distributions

- fitting a distribution: represent data as distribution
- sampling: create new data points that follow the distribution
- factorizing: technique to reduce complexity like joint distributions.

$$P(Y|X) = \frac{P(X|Y) P(Y)}{P(X)} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence.}}$$

- logistic regression fits single regression curve to data
- bayesian " " " distribution of curves.

CDF $P(Z)$ is a running sum that accumulates the probabilities.



drawing a uniform number is always possible.

Sampling Algo

1. discretize the domain of PDF $p(x)$
2. Linearize $p(x)$ if it is multidimensional
3. Calc. CDF $P(z)$ of $p(x)$
4. Draw a uniformly distributed num. u b/w 0 & 1.

$$z^* = \underset{z}{\operatorname{argmin}}_u p(z) \quad u$$

Sampling: picks a random number in the domain of distribution based on the prob. described by it, i.e. picks random number based on prob. described by curve in PDF.

Density Estimation

Creates a PDF from samples. It is opposite of sampling.

parametric density estimation → from pattern recognition.

$$\theta^* = \underset{\theta}{\operatorname{argmax}} p(x_1, \dots, x_n | \theta) \rightarrow \text{Max likelihood estimator.}$$

1. parametric DE req. good func. representation
2. Non-parametric DE can operate on arbitrary distribution.

↳ Ex: Histograms.

Number of bins to choose for histogram is a hyper-parameter.

Non-parametric DE

1. kernel-based

2. nearest neighbors

Let $p(x)$ be a PDF in 0-dim space. ΔR a region around x .
prob. loss in R is $P = \int_R p(x) dx$ → relative freq.

Assu1:

$$P = \frac{\# \text{ points in } R}{\text{total # of points}} = \frac{K}{N}$$

Assu2: R is small

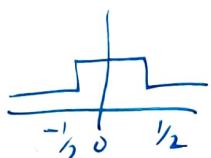
$$P = \int_R p(x) dx = p(x) \int_R dx = p(x) \cdot V \quad \hookrightarrow \text{volume}$$

$$\text{Assu 1+2} \Rightarrow p(x) = \frac{k}{N \cdot V}$$

Parzen Window estimator

Fines V & leaves k/V . [k varies]

$$h(u) = \begin{cases} 1 & \text{if } |u| \leq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$



$$K(x) = \sum_{i=1}^N K\left(\frac{x-x_i}{h}\right)$$

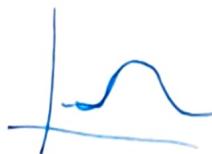
→ scales the box size
= volume of region

$$P(N) = \frac{1}{N} \sum_{i=1}^N \frac{1}{h} K\left(\frac{x-x_i}{h}\right)$$

Kernel removes discretization error of fixed-discrete histogram bins but is still blocky.

gaussian kernel

$$P(x) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2\pi h^2} \cdot \exp\left(\frac{-\|x - x_i\|_2^2}{2h^2}\right)$$



new estimate
will look less bony



Any kernel can be used which satisfies

$$h(u) \geq 0$$

$$\int h(u) = 1$$

K-NN DE

Fixes K & varies V . Calculates V from distance of K -nearest neighbors

Model selection prob

optimizing hyperparameters / model selection

held out part
of training

- hyper-parameters should be tuned on validation set because a training set be most complex model wins as it achieves lowest error.
 - CV(cross validation) is used only for supervised tasks as it requires an objective function for ML (maximum likelihood)
 - Density estimation is unsupervised.
 - Use CV when training data is ltd.
- KNN & parzen window estimators are memory intensive
i.e. need to store all samples.
Also are non-parametric DE.

Bias & variance [Applicable to both supervised & unsupervised learning]

Motivation behind hyperparameter optimization is to aim for generalization to new data.

For kernel DE, pitfalls are

- Too large kernel: covers all space with prob. mass but density is too uniform.
- Too small kernel: closely represents data but might assign too low prob. in areas without training data.
- optimal kernel size represents structure of training data & any unobserved data to some extent.

Bias-Variance trade off → cross validation helps in finding right balance.

Bias → square of average deviation of an estimator from ground truth.

$$\left[\sum_{i=1}^N \frac{(y_i^* - y_i)^2}{N} \right]^2 \rightarrow ?$$

Variance → variance of the estimates.

expected squared deviation from the estimated mean.

high bias → model undercomplexity: doesn't fit to data.

high variance → model overcomplexity: models structure of data & its noise.

Higher model complexity = lower bias & higher variance.

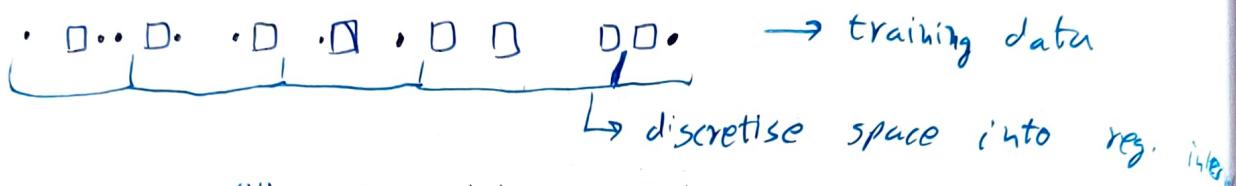
larger kernel support / R increases bias & lowers variance
smaller " " " " lowers bias & increases variance.

Variance & bias require an error metric to be properly understood & quantified. [Ex metrics are MSE & MAE (Mean Absolute Error)]

Ex: unsupervised methods → reconstruction error for PCA
sum or squared dist. from point to μ in k-means clustering.

Random Forests & Variants

decision trees



Histogram → partitions space into equi-distant bins.

$$(3, 2, 2, 0, 1) \quad (1, 2, 1, 3, 1)$$

With Max. likelihood

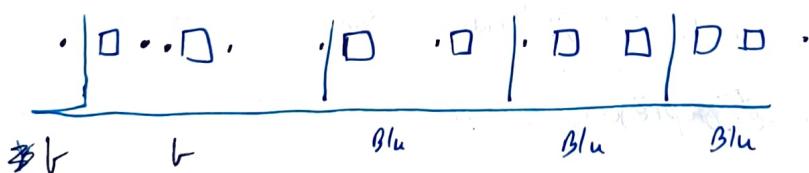
$$\overbrace{(3, 2, 2, 0, 1)}^5 \quad \overbrace{(1, 2, 1, 3, 1)}^5$$

classing bins as $(0, (50-50), 1, 1, (50-50))$

- histogram not flexible. Red & blue appear to form 2 continuous dist, but histogram doesn't represent them fully.

Tree-based representation

Can split the domain into uneven patterns/bins.



Size of bin is derived by an optimisation problem in a data-driven approach.

Forest → number of trees ^{together} is a forest.
It smoothens the boundaries.

- 1. tree-based sample representation
- 2. compare kernels & treebased methods
 - do not need to store data points
 - divide sample space dynamically with objective func.
- 3. CART [Classification & Regression Trees]
 - Random Forests (RF)
 - density forests

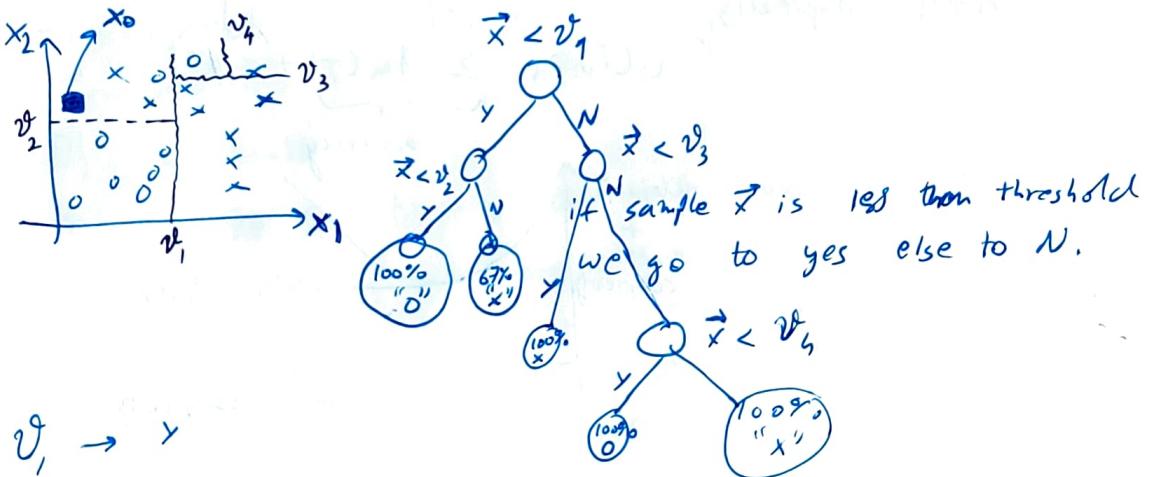
RFs are actively used & not outdated/replaced by Deep neural N/w's.

AIM: Understand CART & RF as tools for hierarchical sample space partitioning.

CART

- recursively divide sample space
- perform classification or regression locally in each partition.

EX



$$x_0 < v_1 \rightarrow Y$$

$$x_0 < v_2 \rightarrow N$$

$x_0 = 67\% "X"$ \Rightarrow classification returns locally observed labels.
similar to K-NN.

Details similar for regression tree.

splits 1-d space into different parts.

const function value as regression estimate.

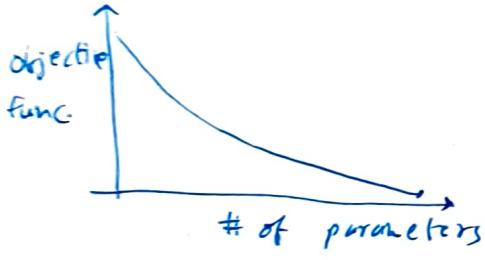
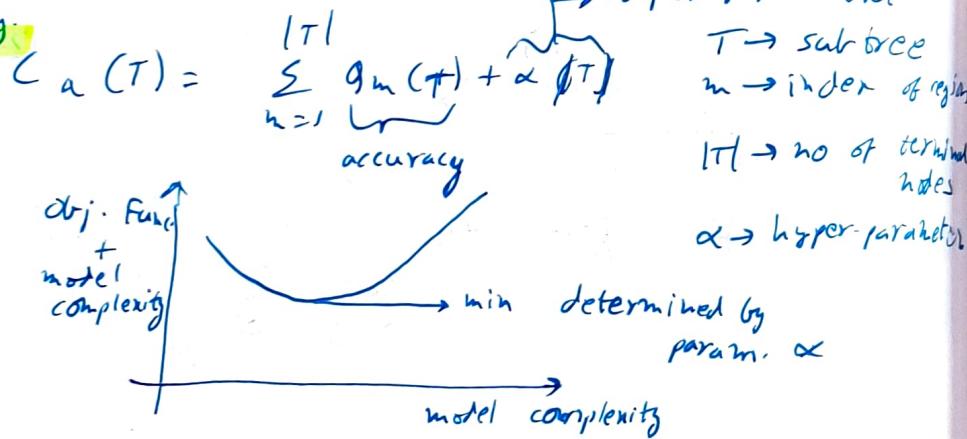
~~random forests~~ fixed value/segment is enough as averaging
multiple trees gives a smooth O/P.

Training

1. each node has 2 parameters
 1. axis-aligned split along one dimension
 2. axis A where threshold from prev. example.
2. tree has max height [Hyper-parameter] → can prevent overfitting.
3. Steps
 - i) grow tree greedily from root to leaves
 - at each node perform exhaustive search for best split w.r.t. on obj. func. $\Omega_m(T)$
 - stop splitting only if few samples are remaining.
 - ii) Prune tree: reduce tree size.

merging splits which do not help solve the task.

pruning cost function is tradeoff b/w model accuracy & model complexity.



Objective functions

Regression

$$\Omega_m(T) = \sum_{x_i \in R_m} (y_i - \bar{y}_m)^2 \quad \bar{y}_m \rightarrow \text{average of all samples in } R_m$$

-ve cross entropy: classification

$$\Omega_m(T) = \sum_k p_{mk} \ln p_{mk} \quad K \rightarrow \text{no. of classes}$$

$p_{mk} \rightarrow \text{relative freq. of class } k \text{ in } R_m$

Classification: Gini index

$$G_2(T) = \sum_k \frac{P_{mk}}{N} (1 - P_{mk})$$

Decision for best splitting candidate $N \rightarrow$ total samples

$$\min \left(\frac{|S_L|}{N} G_L(T) + \frac{|S_R|}{N} G_R(T) \right) \quad l \rightarrow \text{left}, r \rightarrow \text{right}$$

minimises \log_2

searches for candidate split with min. \log_2

$$S_{LR} = \{x_i \in R_{LR}\}$$

Cross entropy & Gini-index are better than misclassification rate as the reward is lower for low purity & high reward for high purity.

Prediction in leaf node

Classification: relative freq.

Regression: average of data points.

RF → more modern than classical CART. \hookrightarrow tendency to overfit.

RF is ensemble of CARTs.

1 tree is weak learner, ensemble is strong.

ensemble averages votes of individual trees. RF, an ensemble of randomized CART to counter overfitting.

Randomising components

- bagging: training on random subset of training data
- random subspace projections in each node
do not do exhaustive search.
only select d-dimensions for best split
- replace exhaustive search with randomly drawn split candidates.

~~Part of RF (frustrate belief in randomisation)~~

What is correlation?

To decrease variance
in RF:

- increase no. of trees; errors caused by individual trees cancel out. Avg. out prediction of more individual models
- use more data! model can learn more accurately
- tune hyper-parameters; prevent overfitting to data

Randomising required
for (trees are less correlated)

- reducing overfitting
- RF are less sensitive to noise
- enhances generalisation

Rationale behind Randomization.

- adds statistical independence & decorrelates trees.

- averaging votes of N independent & iid learners with accuracy $> 50\%$, reduces variance. \rightarrow identical independent

increasing
N no. of trees
reduces variance

$$\sigma_{\text{iid ensemble}}^2 = \frac{1}{N} \sigma^2$$

- iid. is too optimistic. If there is correlation ρ .

$$\sigma_{\text{i.d. ensemble}}^2 = \rho \sigma^2 + \frac{1-\rho}{N} \sigma^2$$

socket variance.

as $N \rightarrow \infty$ $\sigma_{\text{id}}^2 = \rho \sigma^2$ atleast.

Other randomisation steps reduce the correlation ρ .

Configuration of randomisation are additional hyper-parameters:
bagging %, subspace dimension, no. of candidate splits.

Engineering

- cross validation without separate validation set.
use out-of-bag samples for validation

- more trees increase variance & prevent overfitting.

- deeper trees overfit

- parallelisation is easy.

- RF can be used for multiclass classification unlike SVM [space vector machine]
- decision boundaries of RF are smoother. A increase robustness

- trees & forests have explainable decisions.

RF can be seen as data-driven (adaptive) variant to Nearest Neighbour (NN) method.

Density forests

- exchange supervised objective func. with unsupervised function to perform DE.
- unsupervised trees can work only on structure/distribution of data.
- criminisi/shotton/komogorov propose to prefer splits that lead to compact gaussian distributions.
 - compactness is measured using determinant of the gaussian covariance.
 - Re-define **entropy**:
$$G_n(T) = \frac{1}{2} \log ((2\pi e)^n |\Lambda(x_i \in K_n)|)$$

(.1 \rightarrow det
 $\Lambda \rightarrow$ covariance)
 - geometrically, det. approximates rectangular volume around majority samples.
- single tree models the density as multiple truncated trees.
- average over 10 forests smooths out sharp boundaries.
- not necessary to sample from this density.
 1. Randomly select a tree
 2. " move from root to leaf with prob.
$$\frac{|S_L|}{|S_L + S_R|}$$
 where $S_{L/R} = \{x_i \in R_{L/R}\}$
 3. " Sample from covariance in the leaf.
 4. Repeat 3 until the drawn sample is inside of the leaf region.

Training error: error a model makes on the same data it was trained on.
Increasing depth of tree does not increase training error but typically reduces it as the tree can fit better to the data.

- ~~Ridge~~
- Tree based method
- do not need to store ^{all} samples.
 - subdivide sample space dynamically with objective function
 - Only representation but no info reg. dist. of sample.

1. Clustering

1. Variational inference (VI) : to draw conclusion on data
2. appx. technique for inference tasks
3. Markov chain Monte carlo : (MCMC)
4. Manifold learning : reduces dimensionality while preserving structure of data.

Clustering : **Segments** data into meaningful groups

- Assign identical labels to similar samples.
- **Unsupervised** technique
- used for exploring data
 1. K-means
 2. Gaussian mixture model
 3. Mean shift

Manifold learning

- represent data manifold in a lower dim. space.
1. PCA
 2. Multi-dimensional scaling
 3. ISOMAP
 4. Laplacian Eigenmap

K-means Clustering.

- represents observations as a few groups.
- hard clustering method: each sample gets a discrete cluster label.
- popular for: conceptual simplicity
coding simplicity
reasonably fast execution time.
- EX: vector quantization of colors.

K-means will converge even on uniform data although the result would not make sense. Uniform data has least structure possible.

Ideas / objective funct

- cluster is collection of points that are closest to a certain center point.
- assures there are K centers [K-means clustering]
- we can iteratively update each center from the points close to the center.



withincluster distance

$$W(c) = \frac{1}{2} \cdot \sum_{k=1}^K \sum_{q(i)=k} \sum_{q(j)=k} \|x_i - x_j\|^2$$

at sum of euclidean dist of pairs of points within same cluster.

$$= \sum_{k=1}^K N_k \sum_{q(i)=k} \|x_i - \mu_k\|^2$$

$\mu_k \rightarrow$ mean of all points of cluster k.

$N_k \rightarrow$ no. of points

K-Means always converges as with each iteration the $w(c)$ decreases. However it may converge to local minima rather than global minima.

$c(i) \rightarrow$ cluster id

$x_i \rightarrow$ sample in cluster.

locally optimal; K-means is locally optimal

different initializations could result in different results.

The clusters partitions the space: this partitioning is called

Voronoi tessellation

Convergence: change in mean vectors/points changes less than a threshold chosen [ϵ]

GMM

- covered in PR

- Models a PDF as sum of K normal distributions N with weights π_k

$$\text{p}(x) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k) \quad 0 \leq \pi_k \leq 1$$

$$\sum_{k=1}^K \pi_k = 1$$

- can be seen as parametric DE or as a clustering method
- clustering view: one mixture component $\pi_k N(x | \mu_k, \Sigma_k)$ is one cluster.

Model of observations A hidden variables Max Likelihood

$$f(n) = \frac{1}{2\pi\sigma^2} e^{-\frac{(n-\mu)^2}{2\sigma^2}}$$

- direct fitting of GMM parameters via ~~MLE~~ fails due to singularities: likelihood $\rightarrow \infty$ when one component has only 1 observation. [Gaussian to cover all data ~~vs~~ very narrow distribution centered on 1 point]

- instead use Expectation-Maximization (EM) algorithm. [similar to Max Likelihood method]

- hidden variable $p(x) = \sum_z p(x|z)$

introduced by expressing the density as a marginal of a larger joint distribution.

step 1: define GMM parameters in terms of marginal.

step 2: define numerical solver i.e. EM algo

$$p(n, z) = \cancel{x(z)p(n)} p(x|z) p(z) \quad [\text{product rule}]$$

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

$$\sum_z p(x, z) = \sum_z p(z) p(x|z)$$

if z is on or out of n indicator

$$p(z) = \pi \quad (\text{marginalized})$$

z is an indicator that tells us in which component we are in.

$$3: p(x|z) = N(\mu_z, \Sigma_z)$$

- set z as binary indicator vector

it indicates to which component it belongs

$$z_k = \{0, 1\}$$

$$\sum_{k=1}^K z_k = 1$$

z_k selects/reactivates all components except one k or current component

use zero technique

$$p(x|z) = \prod_{k=1}^K N(x|\mu_k, \Sigma_k)^{z_k}$$

resulting in single gaussian at $z_k = 1$

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

- set prior so that its marginal is π_k

i.e. $p(z_k=1) = \pi_k$

$$P(z) = \prod_{k=1}^K \pi_k^{z_k}$$

- Introduce z , then cancel it to make hidden param. explicit, to include them in numerical optimization.

- this method is soft-clustering variant of k-means.

- responsibility indicates degree of membership of a sample to a component. $p(z_k=1|x) = \gamma(z_k)$ responsibility: prob that sample x belongs to component k .

- Annealed sampling is used for sampling from GMM.

sampling from density forest was similar.

$$p(x) = \sum_z p(x|z) = \sum_z \pi_k p(z) = \prod_{k=1}^K \pi_k^{z_k} N(x|\mu_k, \Sigma_k)^{z_k}$$

EM algo:

1. Initialize $\Theta = (\pi_k, \mu_k, \Sigma_k)$

2. Expectation: determine membership of sample to GMM component, i.e posterior distr. of latent variable $p(z|x, \Theta)$

3. Maximization: from those memberships find Θ^{new} via maximizing expectation of complete-data log likelihood $\sum_z p(z|x, \Theta)$

4. set $\Theta = \Theta^{\text{new}}$ & goto 2 until convergence.

Model selection algorithms for GMM such as MCMC A VI
work towards obtaining the posterior of latent variable.

Mean shift algorithm

- iteratively ascends the gradient of a kernel density estimate
- converges to ^{mode}(local max) of density without estimating full density, we want to get to a densest point/cluster centre.
- middle way b/w k-means & GMMs.
- clusters result from a few local kernel DEs.

Kernel notation & constraints

- has to be **radially symmetric** [must depend only on euclidean dist to a sample]
- must accept **squared diff.** $\rightarrow \gamma_p$.

$$K(x_0, x) = c \cdot k(\|x_0 - x\|_2^2)$$

gaussian kernel $k_{\text{gauss}}(x) = e^{-\frac{1}{2}x}$; std dev. is set to ~~1~~ 1

$$k_{\text{epanechnikov}}(x) = \begin{cases} 1-x & \text{for } |x| \leq 1 \\ 0 & \text{else} \end{cases}$$

- kernel sizes: assume $\|x_0 - x\|^2$ are already size **normals**
- smaller kernel give more local clusters
- larger clusters start to merge. Large kernels give less clusters
- do not have to **preset no. of clusters**

KMeans

- result from min dist. b/w samples

GMM

- result from full density estimation

Mean shift

- result from local kernel density estimates.

Using kernel density estimates to go to a local mode of density via gradient ascent.

gradient

$$\nabla p(x) = \nabla \frac{1}{N} \sum_{i=1}^N k_\gamma(x_i, x) = \frac{1}{N} \sum_{i=1}^N \nabla k_\gamma(x_i, x)$$

∇ can be taken in due to linearity.

i.e., gradient can be calculated on the kernel function.

using kernel profile $k_\gamma(x_0, x) = C(\|x_0 - x\|_2^2)$ A substitute $s = \|x_0 - x\|_2^2$

$$\frac{\delta k(s)}{\delta(s)} = k'(s) ; \quad \frac{\delta(s)}{\delta(x)} = \frac{\delta(x_0 - x)^T(x_0 - x)}{\delta x} = -2(x_0 - x)$$

! → defined set to operator

$$\therefore \nabla p(x) = \frac{1}{N} \sum_{i=1}^N C \cdot k'(\|x_i - x\|_2^2) (-2(x_i - x)) = 0$$

$$\sum_{i=1}^N k'(\|x_i - x\|_2^2) (-2(x_i - x)) = 0$$

$$\Rightarrow \sum_{i=1}^N k'(\|x_i - x\|_2^2) x_i - \sum_{i=1}^N k'(\|x_i - x\|_2^2) x = 0$$

Algo:

$$\Rightarrow \frac{\sum_{i=1}^N k'(\|x_i - x\|_2^2) x_i}{\sum_{i=1}^N k'(\|x_i - x\|_2^2)} - x = 0 = m^{(t)}(x)$$

take to other side

- calc mean shift vector
- one gradient ascent step
- update pos $x^{(t)}$

$$x^{(t+1)} = x^{(t)} + m^{(t)}(x) = \frac{\sum_{i=1}^N k'(\|x_i - x^{(t)}\|_2^2) \cdot x_i}{\sum_{i=1}^N k'(\|x_i - x^{(t)}\|_2^2)}$$

- repeat 1 & 2 until convergence.
- [until gradient = 0]

parameters reqd:

- kernel parameters [window size]

- cluster linking param for post processing

[ex: dist threshold]

for grouping data points into clusters

- group samples converging to same location

- with K_{EP} , the update $\pi^{(t+1)}$ is just a calc. of mean.
- Euclidean dist. is sensitive to scaling differences \rightarrow normalize the dimensions of the samples

Model selection for K-means

No objective func. as it is unsupervised.

No labels to compare.

What K to choose?

Gap statistics \rightarrow statistical way to choose K .

Ideas \rightarrow examine WC distance for different K .

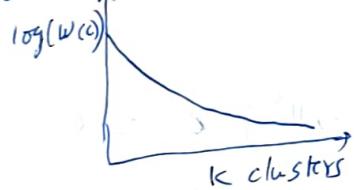
Select smallest K for which $W(c)$ is substantially better

than $W(c)$ of $K+1$ clusters.

natural choice for compact clusters

$$W(c) = \sum_{k=1}^{NK} N_k \sum_{c(i)=k} \|x_i - \mu_k\|^2$$

For increasing k , $W(c)$ has to decrease



- WC cannot search for K
 $W(c)$ as it will be number of samples, i.e. group in supply
 in n clusters which is not sensible. [optimal K cannot be no. of samples]

- elbow method is not good too: searches for bend on $W(c)$ curve.
 it would be unsure which bend is better.

Gap statistics [captures structure in data] [tweaked usage of $W(c)$ to find optimal K]

- relate $W(c)$ of actual samples to $W(c)$ of an artificially created reference.

- reference samples drawn from uniform distribution, representing K clusters [no natural clusters/structure]

- with least possible structure [use 20]

Algo 1. draw B sets of uniform distributions for different K values

2. On those distr. calc. $\log(W(c))$ for different K

denote result $\log(W_{unit}(c))$

3. For the K clusters calc. gap $g(K)$ as diff w.r.t ref. $\log(W_{unit}(c)) + \log(W(c))$.

4. Select K as $K^* = \arg\min_K \{K | g(K) \geq g(K+1) - s_{K+1}\}$

$s_{K+1} = sk \sqrt{1+B} \rightarrow$ unbiased estimate of std. deviation

the actual $\log(w(c))$ is lower than $\log(w_{\text{unit}}(c))$ because
the ref. has least possible structure [unif. dist.]

Relation with Huffman code? think of entropy large (in) uniform
small (in) non-uniform

For the graph choose K where the condition in point 4
is met the first time, because we have to choose smallest

K .

Why choose smallest K ?

- Ockham's razor; choose simplest model unless we have good reason otherwise.

Gap statistics formulates model selection as "finding smallest clustering
that finds notable structure"

Sampling & GMM MCMC inference

choose K , number of components.
complicated than k-means, can be done in a fully probabilistic way.
- GMM fitting with fixed K had no closed form sol but iterative EM algo,

- 1. Sol using additional approximation for selecting K .
- 2. MCMC: approximates intractable func. with finite samples
- 3. Variational Inference: approximates intractable func. with simpler tractable func.

Why we need sampling? Motivation behind Sampling

- TO replace analytical sol.

- if we cannot simulate

then can we approximate by drawing samples?

Eg: expectation $E[f] = \int f(z) p(z) dz$ cannot be solved analytically.

but draw L samples from $p(z)$ & calc.

$$\hat{f} = \frac{1}{L} \sum_{l=1}^L f(z^{(l)})$$

1st time this estimator is unbiased $\Rightarrow E[\hat{f}] = E[f]$ & variance

depends on L which gives best

K for which $w(c)$ is min

so 1st accuracy of estimator

depends on dimension of z ,

hence the gap decreases i.e.

occurrence of minima of $w(c)$ is large

$$\text{var}[\hat{f}] = \frac{1}{L} E[(f - E(f))^2]$$

i.e. few samples suffice

Sampler from Lecture-1 can operate on general distr.
but its usefulness is lth.

- req. full density rep. at every location.
- whatever our density rep. is, we have to convert it to histogram.

Histograms are either

- quite coarse [few bins]
- quite inefficient (B bins); storage is not prob but each bin has to be filled with data points for sufficient statistics.
Ex: 1 billion bins req 10 billion datapts.

Tadoff we make here propagates into approximation error of sampler.

Alternatives:

parametric alternative [parametric std. distr.]

- from uniform distr. to other distr. exist
- Analytic mapping:
 - gaussian
 - exponential
 - cauchy,
 - Algo
 1. sample $p(z)$ from uniform distr.
 2. Transform sample into target distr. $p(y)$ with analytic mapping
 $y = f(z)$

note: When converting one distr. to other, have to take account of the derivative.

$$p(y) = p(z) \left| \frac{dz}{dy} \right|$$

$$y(y_1, \dots, y_m) = p(z_1, \dots, z_m) \left| \frac{s(z_1 - z_m)}{s(y_1, \dots, y_m)} \right|$$

similar to sampler in 1st lecture, instead of calc. CDF calc mapping

$$y = f(z)$$

- Rejection sampling → ie. cannot apply parametric methods
distribution are often more complicated & hence it may be possible
to obtain $p(z)$ upto unknown norm. factor $n(z) = \frac{1}{Z} p(z)$
- define a simpler distr. $q(z) + \text{const } k$ as envelope to $k q(z)$
 - define $\tilde{p}(z)$, s.t. $k q(z) \geq \tilde{p}(z)$ for all z .



$q \rightarrow$ analytic function.

1. draw sample z_0 from $q(z)$
2. draw u_0 from uniform distr. $[0, \log(p_0)]$
3. reject (z_0, u_0) if $u_0 > \tilde{p}(z_0)$ otherwise return z_0
 $u_0 \rightarrow$ quantile value of sample in uniform distr.

certain prob mass of the drawn sample belongs to
 $p(z) \wedge$ some to $\tilde{p}(z)$.

- Why is this method correct?
- prior to rejection: the pair (z_0, u_0) is uniformly distributed across the area of the curve $\log(p)$
 - after rejection, the pair (z_0, u_0) is uniformly distr. across the area of the curve $\tilde{p}(z)$

Adaptive Rejection Sampling
 Rejection sampling gets inefficient if $q \wedge p$ differ too much, ie.
 more samples have to be discarded.
 So this method form a envelope using piece-wise linear functions
 in log space, which is better fitting.

But q might not have a simple analytic form.
 work well on log concave functions, ie. derivatives of $\log(p)$ are
 non increasing function of z .
 fitting lines to log of func. = fitting piecewise exponential
 dist. to orig. func.

$$q(z) = k_i x_i \exp\{-x_i(z - z_i)\} \quad \begin{cases} z_{i-1, l} < z \leq z_{i, i+1} \end{cases}$$

Sampling in Models with Many variables/attributes

Cases:

1. Random variable has many attributes $x = (x_1, \dots, x_n)$ or model consists of many dependent RV. x_1, \dots, x_m
2. corresponds to modelling a sampling based sol. for GMM fitting with model selection: $k, \pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K$ A **hyperprior** are all RVs

Such cases can be called **high-dimensional spaces**.
 One sample is then full set of assignments to all unknowns
 impossible to use rejection sampling + ASR because gap b/w $q \wedge \tilde{p}$ increases

Ex: envelopes around all GMN var.

2. Draw a set of variables A one of them will be likely outside of \tilde{P} . Note variables more likely

MC MC sampling

Mitigates issues in high-dimensional spaces.

Idee:

- sample from 1 var at a time
- repeat this sampling in an iterative manner, using recently sampled values.

- Sample in iteration τ from a state space of variable $z^{(\tau)}$ using prev. iterations $z^{(1)} \dots z^{(\tau-1)}$

- **Markov chain**, models only 1st order statistical dependencies, i.e. weak look at prev. states,

$$\text{skipped } p(z^{(\tau+1)} | z^{(1)} \dots z^{(\tau)}) = p(z^{(\tau+1)} | z^{(\tau)})$$

Ex: Metropolis-Hastings method to consider all past states = considering 1 prev state

Gibbs Sampling

sample distr. $p(z) = p(z_1, \dots, z_n)$ of M R.V. [somehow initialised]

In each step update 1 variable by drawing from that variable conditioned on others.

$$\text{sample } z_1^{(\tau+1)} \sim p(z_1 | z_2^{(\tau)}, \dots, z_n^{(\tau)})$$

$$z_2^{(\tau+1)} \sim p(z_2 | z_1^{(\tau)}, z_3^{(\tau)}, \dots, z_n^{(\tau)})$$

$$z_n^{(\tau+1)} \sim p(z_n | z_1^{(\tau)}, \dots, z_{n-1}^{(\tau)})$$

update for that var
a random draw from dist. of that var. conditioned on all other variables.

→ subsequent samples are correlated due to markov chain.

- after T iterations sample of complex distributions is a single state.
- after $M \cdot T$ sub-draws
i.e. after M \downarrow iterations
No. of RV

Bayesian HMM fitting: Model setup

- Paul HMM parameter π, μ_k, Σ_k rep. a prior distribution
- use conj. prior: ^{param.} distribution \times prior = same distr. family as parameter distr.
- parameters to prior distr. are set to fixed values. [hyperpriors]
- value is not important
- increasing no. of observations overwrites the prior.
- priors shape distribution absence of observations.

Kullback-Leibler Divergence

describes how different 2 distributions $p \neq q$ are

$$KL(p||q) = \sum p(z) \log \left(\frac{p(z)}{q(z)} \right) = - \sum p(z) \log \left(\frac{q(z)}{p(z)} \right)$$

It is the expected log difference w.r.t. p .

derivation: let X be random discrete var with states x_i & $p(x=x_i)=p_i$

$$H(p) = - \sum p_i \log(p_i) \rightarrow \text{lower bound for encoding } X.$$

Without optimal P having encoder q , the relative entropy is:

$$KL(p||q) = \int p(x) \log q(x) dx - \left(\int p(x) \log p(x) dx \right) ; \text{ this is the additional effort required for encoding } p$$

Relative entropy \Leftrightarrow Kullback-Leibler divergence

Expected number of extra bits for encoding text with distribution p , with an optimal code for distribution q .

Properties: 1) $KL(p||p) = 0$

(non-ve) 2) $KL(p||q) \geq 0$; equal only if $p=q$

(non-symmetric) 3) $KL(p||q) \neq KL(q||p)$; it is not symmetric (hence it is not a distance)

Proof for Non-negativity

Jensen's inequality states: $f(tx_1 + (1-t)x_2) \leq t f(x_1) + (1-t) f(x_2)$
generalised form: $f\left(\sum_i t_i x_i\right) \leq \sum_i t_i f(x_i)$

Mapping $\sum p(z) \log \frac{q(z)}{p(z)}$ to $\sum t_i f(x_i)$

hence

$$\log \sum_z p(z) \frac{q(z)}{p(z)} \leq \sum_z p(z) \log \frac{q(z)}{p(z)}$$

log's concavity to convexity

$\Rightarrow 0$ as summation of $q(z)=1$ & $\log(1)=0$

\therefore KL divergence ≥ 0 always.

Non-Symmetry

- When fitting a distribution $q(x)$ to a fixed distribution $p(x)$ by minimising KL divergence, $KL(q||p)$: it focuses ~~on~~ q to a more of p , to avoid probability mass of q anywhere where p is close to 0.
- Minimizing reverse $KL(p||q)$ leads to a q that is broadly covering p , i.e. q distributes prob. mass everywhere where ~~is~~ p is non-zero.

General Expectation Maximisation

It is a solver for models consisting

- observations, X
- hidden variables, Z
- parameters θ

When it is not possible to fit θ to X , i.e. difficult/infeasible to solve directly for $p(X|\theta)$.

In general we seek to maximise likelihood $p(X|\theta)$ as

$$p(X|\theta) = \sum_z p(X, z|\theta)$$

where z is hidden variables to simplify fit by splitting up the 2 steps

1. E step: improves z via $p(z|x, \theta)$
2. M step: improves θ via

Decomposition: log evidence decomposes to ELBO & KL divergence

$$\log p(x|\theta) = \log \sum_z p(x, z|\theta) = L(q, \theta) + KL(q||p)$$

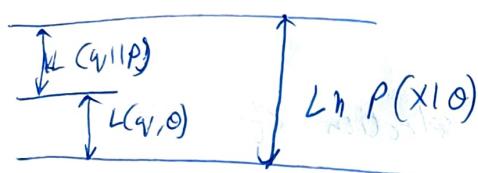
achieved after introducing a distribution over latent variables $q(z)$

where:

$$L(q, \theta) = \sum_z q(z) \log \left\{ \frac{p(x, z|\theta)}{q(z)} \right\}$$

$$KL(q||p) = - \sum_z q(z) \log \left\{ \frac{p(z|x, \theta)}{q(z)} \right\}$$

to undo decomposition: put $\log p(x, z|\theta) = \log p(z|x, \theta) + \log p(x|\theta)$ in ELBO



goal is to iteratively improve $q(z)$ & θ in $L(q, \theta)$

- steps
1. E step: updates $q(z)$ with θ^{old} , effectively setting $KL(q||p) = 0$, by setting $q(z) = p(z|x, \theta^{\text{old}})$
 2. M step: update θ with $q(z)$ which recreates $KL(q||p) > 0$ at new location $p(x, z|\theta^{\text{new}})$

M-step provides a set of parameters

E-step " a $q(z)$ that closes gap w/ $p(x|\theta)$

Variational Inference

Expresses posterior $p(z|x)$ with a suitably selected & parametrized distribution q .

Used when estimation of original posterior $p(z|x, \theta)$ is infeasible then fit a feasible distribution q .

2 choices for simplification

1. parametric distribution $q_j(z|w)$ where fitting parameter w is tractable
2. factorized distr. $q(z_1, z_2, \dots, z_n) = \prod_{i=1}^n q_i(z_i)$ where fitting independant factors is tractable.

Example for parametric distribution: Entropy

$$H[p] = - \sum_i p(x_i) \log p(x_i)$$

calculus of variations can be used to find $p(x)$ that maximises/minimises $H[p]$

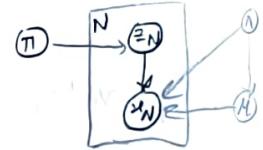
A simpler distribution q could be a Gaussian, such that
min. entropy is reached for $\sigma \rightarrow 0$
max. " " " " " " $\sigma \rightarrow \infty$

Note: q only approximates these solutions &
only models a tiny subspace for all possible distribution q .

For factorized distribution: GMM

- GMM model is complex for automatic selection of no. of components.
- estimating $P(z|x, \theta)$ [posterior] is hard, so std. EM doesn't work
- VI solution approximates original model p with a dist. q that can be factorised into multiple simpler terms.

Full GMM model



$$p(x, z, \pi, \mu, \Lambda) = p(x|z, \mu, \Lambda) p(z|\pi) p(\pi) p(\mu|\Lambda) p(\Lambda)$$

$z \rightarrow$ latent variable, indicates which component creates which sample
 z depends on distribution π [the no. of weight of the components]
component means μ depend on associated standard deviation (Λ)
observations x depend on z, μ, Λ

Factorized approximation
unknown variables captured in variational distribution q

tractable if q can be factorised b/w $z \neq \pi, \mu, \Lambda$
i.e. $q(z, \pi, \mu, \Lambda) = q(z) q(\pi) q(\mu) q(\Lambda)$

$$q(z) q(\pi, \mu, \Lambda) = q(z) q(\pi) q(\mu, \Lambda) [\because \pi, (\mu, \Lambda) are conditionally independent in p]$$

Mean Field approximation

Generalizes EM framework to do variational inference on factorized distribution.

$$q_{\theta}(z) = \prod_{i=1}^M q_i(z_i) \quad \left[\begin{array}{l} \text{factorises all hidden variables into} \\ M \text{ partitions} \end{array} \right]$$

If q_i is simple then $E(\text{expectation})$ can be calculated & reuse the EM framework.

→ IN EM FRAMEWORK

$$L(q_{\theta}, p) = -KL(q_{\theta} || \tilde{p}) - \sum_{V_i \neq V_j} \int q_{\theta}(z) \ln q_{\theta}(z) dz$$

where $\log \tilde{p}(x, z) = \sum_{V_i \neq V_j} \log p(x, z) = \int \log p(x, z) \prod_{V_i \neq V_j} q_{\theta}(z_i) dz_i$

Curse of Dimensionality

Looked at low dimensional feature vectors [also to visualise]
 Real data could consist more dimensions.
 Difficulty increases with data dimensionality.
 ↳ = curse of dimensionality

3 difficulties:

1) Visualization:

- data often more larger than 3D
- To understand data, it helps to visualise/simplify it.

Ex: Remote sensing processes satellite recordings
 [photograph of surface of earth done in several narrow color bands]

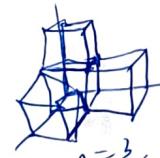
DNN learn featuremaps [representations] with many dimensions

Recommender systems compress large sparse matrices into few product proposals

2) Exponential no. of bins

Assumption: similar features at similar locations in sample space.

∴ classifier or regressor must make local predictions.
 For simplicity equally sized cells: no. grows exponentially with dimensions.



which requires more model parameters & most importantly growing

n. of data points for sufficient observations per cell.
 (Kernel estimators have difficulties in higher dimensions)

3) Distances become less discriminative

- most samples lie at boundary (dist. b/w samples become similar hence dist. are less meaningful in high dimensions)
- rejection sampling is also affected due to this issue

Example/ illustration

D-dim sphere with radius 1, with uniformly distributed samples.

$V_D(r) = k_D \cdot r^D \rightarrow$ general form of volume in D-dim

\nwarrow const. vol. factor

Consider another sphere with radius $1-\epsilon$.

(fraction of spheres is 2)

$$f_D(\epsilon) = \frac{V_D(1) - V_D(1-\epsilon)}{V_D(1)} = \frac{1 - (1-\epsilon)^D}{1} = 1 - (1-\epsilon)^D$$



looking at $f_D(\epsilon)$ it rapidly reaches 1

$$D=10 \quad \epsilon=0.1 \quad f_{10}(0.1) = 65\%$$

$$D=100 \quad \epsilon=0.01 \quad f_{100}(0.01) = 63\%$$

Dimensionality reduction

1. PCA
2. MDS
3. ISOMAP
4. LE

PCA

Also known as Karhunen-Loeve Transform (KLT)

[comes in IVC] Transform coding.

- provides compact representation in lower space.

- PCA is linear projection onto orthogonal basis U

$$U_i^T U_j = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{Otherwise} \end{cases}$$

- basis are eigen vectors of data covariance (Mean-free)
- Idea: normalize data & to perform eigen value decomposition.
- Magnitude of eigen values indicates contribution of a dimension to covariance of data.

Objective function

vectors $x \in \mathbb{R}^d$ d-dimen. space
Find linear mapping $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ $d' \ll d$ [maximises variance of data along each dimension]

\rightarrow obj. func.

$$J = \sum_{i,j=1}^N (\phi_{x_i} - \phi_{x_j})^T (\phi_{x_i} - \phi_{x_j}) + \gamma (\phi^T \phi - 1) ; x_i, x_j \rightarrow \text{data points.}$$

zero-mean samples $\Rightarrow \sum_{i=1}^N x_i = 0$ [in practice subtract mean of all samples]

derivation:

- seek projection of U onto 1D subspace that maximises the variance & show U is largest eigenvector of covariance matrix.

- U is 1st col in ϕ , further vectors iteratively obtained.

- proj data on to $d-1$ dimensional subspace orthogonal to U
- Repeat: \circlearrowleft

To begin $U \in \mathbb{R}^d$ be an arbitrary direction with unit length
i.e. $U^T U = 1$

Inner prod. $U^T x$ projects x onto a 1-D space

Variance of proj data : $\frac{1}{N} \sum_{i=1}^N (U^T x_i - \bar{U} \bar{x})^2 = \bar{U}^T S \bar{U}$; \bar{x} : component wise mean of all x_i
 S : covariance matrix

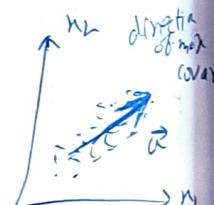
$$S = \frac{1}{N} \sum_{i=1}^N (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})^T$$

to maximise variance of proj data

$$\bar{U}^T S \bar{U} + \lambda \underbrace{(1-U^T U)}_{\text{constraint}} \rightarrow \text{maximise}$$

\rightarrow Lagrangian multiplier to includ constraint $U^T U = 1$,

$$\frac{\delta}{\delta U} U^T S U + \lambda(1-U^T U) \stackrel{!}{=} 0; \text{ for maxima}$$



$$2S\bar{U} = 2\lambda\bar{U}$$

$S\bar{U} = \lambda\bar{U} \Rightarrow$ eigen vector decomposition of S .

i. largest eigen vector associated with largest eigen value gives max covariance. This vector is called principle component.

Approximation error / amount of info lost when operating in low dim. space = $\frac{\sum_{i=1}^d \lambda_i}{\sum_{i=1}^d \lambda_i}$

if = 0.98, then it means subspace preserves 98% of variance of data.

MDS

- equivalent to PCA but operates on distances b/w samples
 - useful if only relative info about samples is available.
(E.g. from a learned similarity measure)
- operates on dist. matrix $D = [d_{ij}^2]$ where $1 \leq i, j \leq N$
where $d_{ij}^2 = (x_i - x_j)^T (x_i - x_j)$; but $x_i \wedge x_j$ are unknown; $x_i, x_j \in \mathbb{R}^d$
- Aim: construct $\tilde{x} = (\tilde{x}_1, \dots, \tilde{x}_N) \in \mathbb{R}^{N \times d'}$ where $d' \ll d$ is an orthogonal proj. onto d' -dim subspace.
- (\tilde{x}_i is unique solution)
 x must be mean free as we cannot distinguish spatial translations. (mean=0).

Matrix Notation for d_{ij}^2 : See notes on MDS

$$d_{ij}^2 = (x_i - x_j)^T (x_i - x_j)$$

$$= x_i^T x_i - x_i^T x_j - x_j^T x_i + x_j^T x_j$$

$$= x_i^T x_i + x_j^T x_j - 2 x_i^T x_j$$

$$D^2 = \text{diag}(x^T x) \cdot \mathbf{1} \mathbf{1}^T + \mathbf{1} \mathbf{1}^T \cdot \text{diag}(x^T x)^T - 2 x^T x \rightarrow \star$$

$$\mathbf{1} = (1, 1, \dots, 1)^T \in \mathbb{R}^N$$

$\text{diag}(x^T x)$ is matrix with principle diag elements of $x^T x$
0 everywhere else.

\star links distances to actual coordinates, i.e. x

recover x by multiplying D^2 with centring matrix C

$$C = \left(I - \frac{1}{N} \mathbf{1} \mathbf{1}^T \right) \quad I \rightarrow N \times N \text{ identity matrix}$$

$$-\frac{1}{2} C D^2 C = -\frac{1}{2} \left(I - \frac{1}{N} \mathbf{1} \mathbf{1}^T \right) \left[\text{diag}(x^T x) \cdot \mathbf{1} \mathbf{1}^T + \mathbf{1} \mathbf{1}^T \cdot \text{diag}(x^T x)^T - 2 x^T x \right] \left(I - \frac{1}{N} \mathbf{1} \mathbf{1}^T \right)$$

simplifies to $\boxed{-\frac{1}{2} C D^2 C = x^T x}$

$\therefore x$ can be obtained from Singular Value Decomposition [SVD]

$$\text{SVD} \left(-\frac{1}{2} C D^2 C \right) = \text{SVD} (x^T x) = U^T \Sigma V \Rightarrow \boxed{x = \Sigma \frac{1}{2} V}$$

$\Sigma \rightarrow$ diag. matrix : $\Sigma^{1/2}$ = square root of diag. entries.

- like in PCA, now select no. of dimension d' from magnitudes of the singular values.
- orthogonal projection to $d' < d$; components can be done by setting all singular values in Σ beyond the first d' -dimensions to 0.

ISOMAP

PCA & MDS perform
Linear mapping.

- performs non-linear mapping.
- Think of swiss-roll [unrolling the swiss roll]
- PCA & MDS squashes the swiss roll.
- Non Linear Mapping approaches use the MDS concept of dist. b/w samples

ISOMAP uses "Non-linearity hack" for MDS.

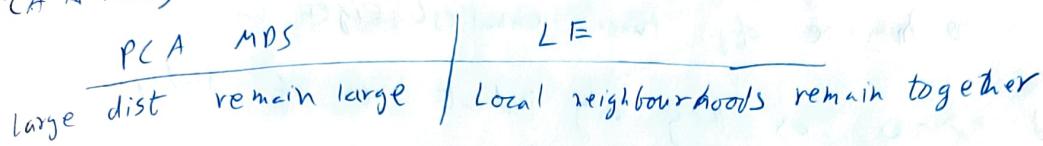
- Preserves non-linear manifold by replacing euclidean dist. with geodesic distances on manifold.
- Dist. is calculated as shortest path from a graph where each sample is connected to its neighbours.

Algorithm:

1. define graph with edge weights as euclidean dist.
(using K-NN or fixed dist. threshold)
2. Distance matrix = all pairs-shortest paths b/w samples (Floyd Warshall algorithm / Minkowski metric)
3. Perform MDS on distance matrix.

LE (Laplacian Eigen Map)

- uses local distances (advances one step from using local distances)
- preserves local neighbourhoods instead of preserving global variance (like PCA & MDS)



- organises local neighbourhoods in a graph Laplacian (general purpose matrix representing graphs)
- more robust to noise than ISOMAP.

Algorithm:

1. Build adjacency graph (via K-NN or fixed dist. threshold)
2. weight the edges using sample similarity. (ex: using kernel) can be replaced by learned edge weights
3. Similarities are called affinities.
4. Perform eigen decomposition of graph Laplacian.
5. Project sample to lower space.
6. Project sample to lower space.

Learn edge weights using density forest

1. train density forest
2. use sample affinity, ex: relative freq. that both samples end up in same leaf node or trees.

Learned weights are advantages as similarity measure better adapts to data.

objective function

$$\sum_{i=1}^N \sum_{j=1}^N (x'_i - x'_j)^2 w_{ij} \rightarrow \text{minimise}$$

$$x'_i, x'_j \in \mathbb{R}^{d'}$$

$$w_{ij} = \exp\left(-\frac{\|x_i - x_j\|_2^2}{2}\right); x_i, x_j \in \mathbb{R}^d$$

solution where $x'_1 = \dots = x'_N$ is prevented with additional cond.

$$\tilde{x}^T D \tilde{x} = 1$$

$$\text{where } \tilde{x} = (x'_1, x'_2, \dots, x'_N)^T$$

$$D = \text{diag}\left(\sum_{i=1}^N w_{1i}, \dots, \sum_{i=1}^N w_{Ni}\right)$$

rewriting obj function

$$\sum_{i=1}^N \sum_{j=1}^N (x'_i - x'_j)^2 w_{ij} = 2 \cdot (\tilde{x}^T (D - W) \tilde{x})$$

$L = D - W \rightarrow$ one variant of graph Laplacian

min. $\tilde{x}^T L \tilde{x}$ subject to $\tilde{x}^T D \tilde{x} = 1$

$$\sum_{i=1}^n (\tilde{x}^T L \tilde{x} + \lambda \tilde{x}^T D \tilde{x}) = 0 \Rightarrow D^{-1} L \tilde{x} = -\lambda \tilde{x}$$

∴ goal is to minimise obj. func. smallest eigen value indicates sd.

- discard $\lambda < 0$ as neg only indicate no. of independent components.

Manifold forests + Laplacian Eigenmaps

- L can be composed of any set of pairwise affinities.

- L can be composed of any set of pairwise affinities defined on a density forest.

Manifold forests have affinities defined on a density forest.
For each tree, define an affinity matrix W_t using dist. $d_t(x_i, x_j)$:

$$W_t = [w_{ij}]_t = \exp(-d_t(x_i, x_j))$$

Gaussian Affinity

$$d_t(x_i, x_j) = \begin{cases} \frac{(x_i - x_j)^T (x_i - x_j)}{\sigma^2}, & \text{if } \text{leaf}(x_i) = \text{leaf}(x_j) \\ \infty, & \text{o.w.} \end{cases}$$

$$\text{Binary Affinity} \\ d_t(x_i, x_j) = \begin{cases} 0, & \text{if } \text{leaf}(x_i) = \text{leaf}(x_j) \\ \infty, & \text{o.w.} \end{cases}$$

with binary affinities, nodes are permuted s.t. W_t is a block matrix

$$W_t = \begin{pmatrix} 1 & 1 & 1 & & & \\ 1 & 1 & 1 & & & \\ 1 & 1 & 1 & & & \\ & & & 1 & 1 & 1 \\ & & & 1 & 1 & 1 \\ & & & 1 & 1 & 1 \\ & & & 1 & 1 & 1 \\ & & & 1 & 1 & 1 \\ & & & 1 & 1 & 1 \\ & & & 1 & 1 & 1 \end{pmatrix}$$

- Averaging $W = \frac{1}{t} \sum_{t=1}^T W_t$

over Forest removes block structure.

- Affinities are large for sample pairs within same leaf. creates sense of locality which is learned from data.

- W can be used to construct L .

$$L = I - D^{-1/2} W D^{-1/2}$$

obtain embedding to d' dimensions by:

1. sort eigenvalues & vectors in increasing order
2. Discard all eigen vectors associated with eigen values $\lambda = 0$
3. lower-dim. coordinates ~~of~~ of x_i are i-th entries of eigen-vectors associated with d' smallest eigenvalues.

Take aways

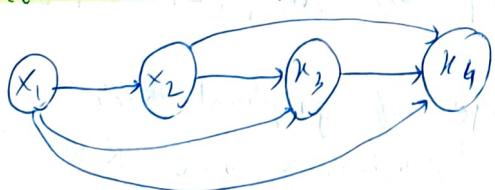
- LE more robust to outliers than ISOMAP
- ISOMAP: shortcut on manifold can be created by few outliers.
- ISOMAP: reduces shortest path.
- Graph Laplacian connects graphs with linear algebra. [tool in spectral graph theory]
- Graph Laplacian connects graphs with linear algebra. [tool in spectral graph theory]
- Ex: K-eigen vectors of largest eigen values of L partition to sample neighbourhood graph. Then K-means assigns clusters.

Intro to Probabilistic Graphical Models

- look at relationship b/w random variables here
 - probabilistic models represent data properties in probabilistic formulations
 - start with joint distribution $p(x_1 \dots x_N)$
- Factorize using product rule
- $$p(x_1 \dots x_N) = p(x_N | x_1 \dots x_{N-1}) \dots p(x_2 | x_1) p(x_1)$$
- order can be arbitrary
hence even this is correct $\Rightarrow p(x_1 \dots x_N) = p(x_3 | x_1, x_2, x_4 \dots x_N) \dots p(x_2 | x_1) p(x_1)$

Directed edges = conditional prob.

Undirected \equiv joint prob.



Ex:

x_1 conditions all variables \rightarrow 3 outgoing edges

x_4 depends on all other variables
 \rightarrow 3 incoming edges

x_2 depends on x_1 A conditions

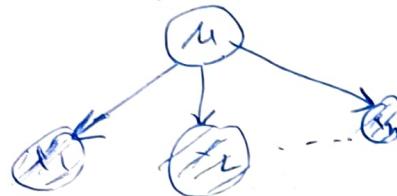
$x_3 \wedge x_4 \rightarrow$ 1 incoming edge +
2 outgoing edges

Example 1 Parametric DE

Gaussian data with fixed σ^2 but unknown μ .

$\mu \rightarrow$ Random Variable

$x_1, \dots, x_N \rightarrow$ mutually independent observations



observed nodes are shaded.

Bayesian DE reverses the arrows to find $p(\mu | x_1, \dots, x_N)$ from likelihood $p(x_1, \dots, x_N | \mu)$. A prior $p(\mu)$.

Plate notation simplifies above graph

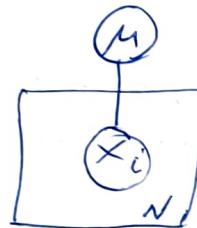
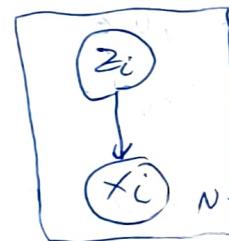


plate is short-hand notation for
N independent, identically distributed
(i.i.d.) of some variables.

Ex 2: Mixtures

- it is an overlay of multiple distributions.

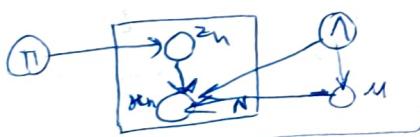
$$p(x_i) = \sum_{z_i} p(x_i | z_i) p(z_i) \quad z_i \rightarrow \text{indicator variable}$$



$p(x_i | z_i = k)$ generates x_i from k-th component.

Ex 3: GMM with hyperpriors

$$p(x, z, \pi, \mu, \Lambda) = p(x | z, \mu, \Lambda) p(z | \pi) p(\pi) p(\mu | \Lambda) p(\Lambda)$$



- Advantage of graphical model
- makes dependencies b/w R.V. more explicit through edges.
- variable independence tractability is better

Conditional Independence

Independence enables factorization of models with many variables.

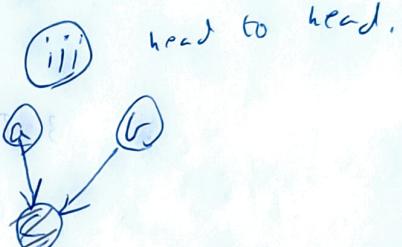
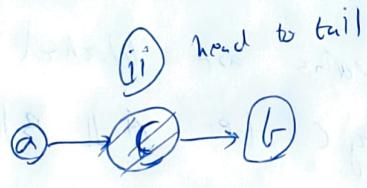
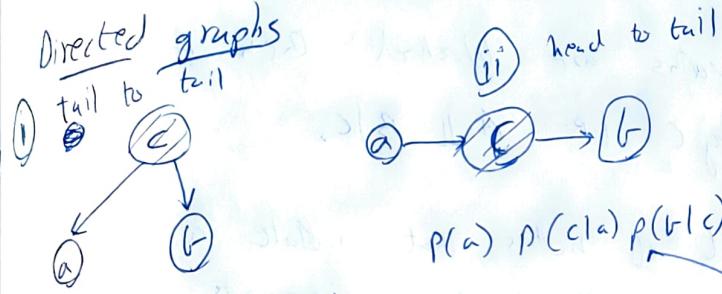
If $p(a|b, c) = p(a|c)$ then ~~a~~ a is conditionally independent of b given c.

$a \perp\!\!\!\perp b | c$ mathematical notation

Tricky on directed graphs / simpler on undirected graphs

Undirected graphs: Markov blanket enables inference.

Directed graphs



$$P(a) P(c|a) P(b|c)$$

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

- if c is observed a, b become independent

$P(a) P(b) P(c|a, b)$
a, b become
(conditionally dependent
if c is observed)

) tail to tail

if c is unobserved ~~$a \perp\!\!\!\perp b | c$~~ since $P(a, b) \neq P(a) P(b)$

$$P(a, b) = \sum_c P(c) P(a|c) P(b|c) \neq P(a) P(b)$$

if c is observed

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

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

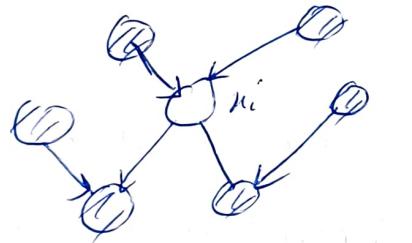
D-separation of var. in directed graph:

- D-sep indicates weaker $\perp\!\!\!\perp$ for sets of A, B, C ; $A \perp\!\!\!\perp B | C$
- vars. in C are observed.

Algo:

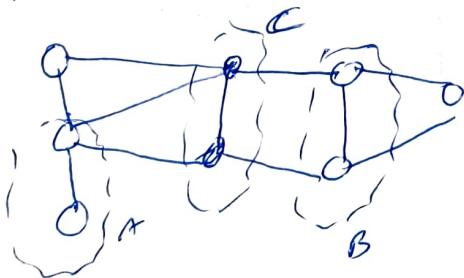
1. consider all path b/w $A \times B$
2. A path is blocked if
 - node is in C & its arrows are ~~head to tail~~ tail to tail or head to tail.
 - node & its descendants are not in C , its arrows are head to head.
3. If all paths are blocked then $A \perp\!\!\!\perp B | C$.

Markov Blanket: minimal set of nodes that isolate a node from rest of the graph. Include immediate parents, immediate children & co-parents of intermediate children (children of children).

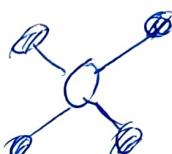


Undirected graph:

separate $A \perp\!\!\!\perp B$ by observed nodes for cond. independence



Markov Blanket: includes set of neighbours of a node,



Hidden Markov Models (HMM)

Any distribution can be represented as a GM. But not practical.

Multiple linked LRU may be more efficient.

Multiple \rightarrow many \rightarrow classical tool to process sequential data.

- classical tool to process sequences
 - can be seen as probabilistic state machine which maps observation sequences to states.
 - states are hidden as user sees only input op. but not states

HMMs are generative probabilistic models that generate samples from trained model.

new samples
Preliminary considerations on seq. data

$$\begin{aligned}
 \text{data} \\
 p(x_1, x_2, \dots, x_T) &= p(x_1) \cdot \prod_{t=2}^T p(x_t | x_1, \dots, x_{t-1}) \\
 &\stackrel{!}{=} p(x_1) \prod_{t=2}^T p(x_t | x_{t-1})
 \end{aligned}$$

1st order Markov Model



- inference can be done easily.
 - ex: if $x_1 \dots x_T$ are measurements on diff. days then treating each x_t (meas.) separately doesn't make sense. Hence introduce dependency on the measurement of ~~prev.~~ day.
convergence saturation quickly.

such model reaches performance level
prev. 2 day's data.

$$p(x_1 \dots x_T) = p(x_1) p(x_2 | x_1) \prod_{i=3}^T p(x_i | x_{i-1}, \dots, x_2)$$



This approach cannot be scaled. Let's keep no. of parameters tractable:

HMM use a trick to keep no. of start variables. (z_t)

use a trick (e.g. never)
they use latent variables. (zt)
→ encode discrete states [April weather
Indian summer]

Match seq. of days with particular weather to plausible seq. of states.

HMM joint distribution

$$P(x_1, \dots, x_T, z_1, \dots, z_T) = P(\underbrace{x_1, \dots, x_T}_{\text{obs. seq.}} | \underbrace{z_1, \dots, z_T}_{\text{hidden state seq.}}) r(z_1, \dots, z_T)$$

Assumptions:

1. each obs. depends on single hidden state

$$P(x_1, \dots, x_T | z_1, \dots, z_T) = \prod_{t=1}^T P(x_t | z_t) \rightarrow ①$$

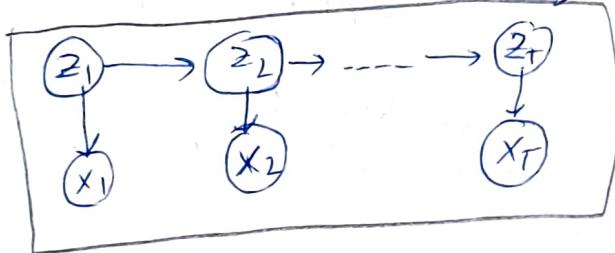
2. each state depends on predecessor

$$P(z_1, \dots, z_T) = P(z_1) \prod_{t=2}^T P(z_t | z_{t-1}) \rightarrow ②$$

From ① using ① & ②

$$\Rightarrow P(x_1, \dots, x_T, z_1, \dots, z_T) = P(z_1) \prod_{t=1}^T P(x_t | z_t) \prod_{t=2}^T P(z_t | z_{t-1})$$

↳ factorized HMM consists of small terms.



can describe complex data in a compact representation via

1. distr. of starting state $P(z_1)$

2. likelihood of state transitions $P(z_t | z_{t-1})$

3. likelihood of observation in each state $P(x_t | z_t)$.

Rabiner Notation

N states s_i

sequence of T observations; each observation is discrete symbol o_v
from alphabet of size V .

HMM parameters $\lambda = (A, B, \pi)$

$A \rightarrow$ state transitions matrix from state s_i to s_j , $A = [a_{s_i s_j}] \in R^{N \times N}$

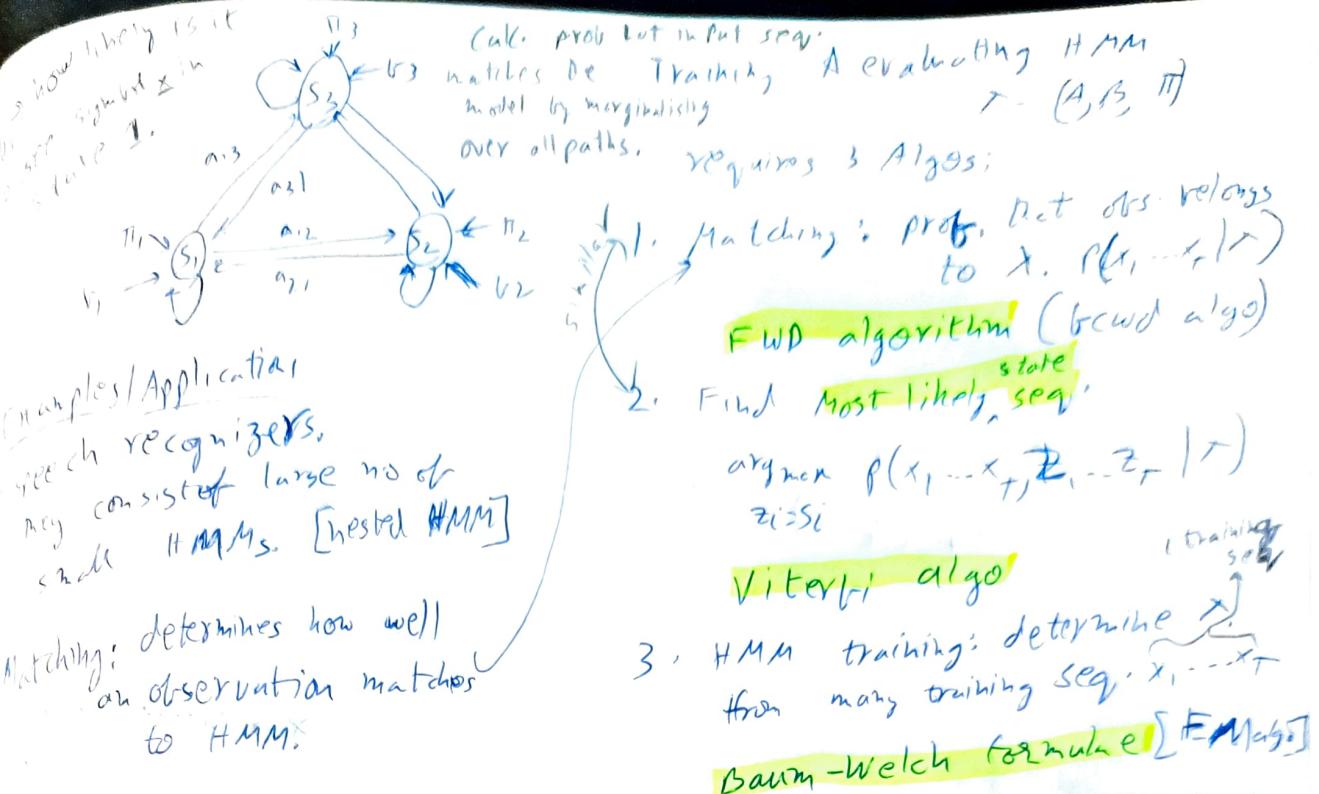
$$\sum_{j=1}^N a_{s_i s_j} = 1$$

$B \rightarrow$ matrix of symbol prob. per state s_i with $\sum_{v=1}^V b_{s_i}(o_v) = 1$

$$B = [b_{s_i}(o_v)] \in R^{N \times V}$$
 matrix of observation we can make

$\pi \rightarrow \pi = [\pi_{s_i}] \in R^N$ vector of prob. that seq. starts in state i .

$$\sum_i \pi_{s_i} = 1$$



conditional dependence of Algo 1 & 2 implies that they work on a trained model.

Algo 1: How well Matches an I/p seq' to specific HMM.

calc. $p(x_1 \dots x_T | \lambda)$ by marginalising over all state seq. $z_1 \dots z_T$

$$p(x_1 \dots x_T) = \sum_{z_1=S_1}^{S_N} \sum_{z_2=S_1}^{S_N} \dots \sum_{z_T=S_1}^{S_N} \prod_{i=1}^T b_{z_i}(x_i) a_{z_i z_{i+1}} b_{z_{i+1}}(x_{i+1}) \dots a_{z_{T-1} z_T} b_{z_T}(x_T)$$

any state can be 1st, 2nd or last state respectively

\Rightarrow computational cost of $O(N^T)$.

state seq' an input takes can have many possibilities.
prob. of 1st obs. in z_1 prob. of making 2nd obs. in z_2

$p(z_1 \dots z_T) = \sum_{z_1=S_1}^{S_N} \sum_{z_2=S_1}^{S_N} \dots \sum_{z_T=S_1}^{S_N} \prod_{i=1}^T b_{z_i}(x_i) a_{z_i z_{i+1}} b_{z_{i+1}}(x_{i+1}) \dots a_{z_{T-1} z_T} b_{z_T}(x_T)$

prob. of starting in state z_1 prob. to start in state z_1 next state could be any other state
any state could be starting state prob. of Trans. from z_1 to z_2

factoring in.

$p(x_1 \dots x_T) = \sum_{z_1=S_1}^{S_N} \prod_{i=1}^T b_{z_i}(x_i) a_{z_i z_{i+1}} b_{z_{i+1}}(x_{i+1}) \dots a_{z_{T-1} z_T} b_{z_T}(x_T)$

poss. to start in state 1 $\sum_{z_1=S_1}^{S_N}$ $\prod_{i=1}^T b_{z_i}(x_i)$ $\sum_{z_2=S_1}^{S_N} a_{z_1 z_2} b_{z_2}(x_2)$ $\dots \sum_{z_T=S_1}^{S_N} a_{z_{T-1} z_T} b_{z_T}(x_T)$

look for assignment to 2nd state $O(N^2 T)$

Algorithm:

1. initialise

$$t=1 \quad \alpha_t(z_1 = s_i) = T_{21} b_{z_1}(x_i)$$

↓
prob we start in state i
match symbol to state i

$$\alpha_1(z_1 = s_1) = 0.1$$

$$\alpha_1(z_1 = s_2) = 0.05$$

$$\alpha_1(z_1 = s_N) = 0.01$$

↓ read from list

2. Iteration: go from $T-1$ to t

$$2 \leq t \leq T; \alpha_t(z_t = s_i) = \sum_{z_{t-1}=s_j}^{s_N} \alpha_{t-1}(z_{t-1}) \cdot a_{z_{t-1}, z_t}$$

at time t we have
prob that we are in
state s_i
we could have come
from any other
states.

↓
prob that in previous time
we were in previous state
read from list

→ caching variable
for path probability

3. summation over end points

$$t=T \quad p(x_1, \dots, x_T) = \sum_{z_T=s_i}^{s_N} \alpha_T(z_t)$$

sum over all possible paths.

Use this algo to see how likely the seqn has been generated by the model.

Also for RE training.

Backward algo: start at time T & goto 1.

1. init \rightarrow caching var. β

$$t=T; \beta_T(z_T = s_i) = 1$$

prob at time T we are in state s_i .

2. Iterate

$$1 \leq t \leq T \quad \beta_t(z_t = s_i) = \sum_{z_{t+1}=s_j}^{s_N} \beta_{t+1}(z_{t+1}) b_{z_{t+1}}(x_{t+1}) a_{z_t, z_{t+1}}$$

↓
at $t+1$ we were in state i

↓
prob for current state to next state

BCKWD algo doesn't use state x_i [OK for training use]

time: $t-1$

state

(s_1)

$\alpha_{t-1}(s_1)$

(s_2)

$\alpha_{t-1}(s_2)$

(\vdots)

$\alpha_{t-1}(s_N)$

prob that we were in state i at $t-1$

$a_{s_1 s_1}$ ↓ we were at s_1 at $t-1$

$a_{s_2 s_1}$ ↓ we were at s_2 at $t-1$

$a_{s_N s_1}$

step: $t-1 \rightarrow t$

$\alpha_t(s_i) = \sum_{z_{t-1}=s_j}^{s_N} \alpha_{t-1}(z_{t-1}) \cdot a_{z_{t-1}, s_i}$

prob at time t we are in state s_i

$\alpha_t(s_1)$ is summed up

repeated for all states

Alg 2: Searches for most likely state seq.

Find: $\arg \max_{z_t = s_j} P(x_1, \dots, x_T, z_1, \dots, z_T | \lambda) \rightarrow$ from prob. of seq. of observation A seq. of states give trained model we seek the assignment of states to hidden variables that maximise prob. [most likely state seq.]

$\delta_t(z_i) \rightarrow$ caches prob. of most likely path to z_i within time 1.

$\psi_t(z_i) \rightarrow$ caches predecessor state to reconstruct path.

Viterbi alg 0:

i. Init : $t=1: \delta_1(z_1 = s_i) = \pi_{z_1} b_{z_1}(x_1)$

and state can be our starting state

starting prob.

$\psi_1(z_1 = s_i) = 0 \rightarrow$ no prev. state prob. of ~~most~~ most likely path that leads us in time to state i

ii. Iteration

$1 \leq t \leq T: \delta_t(z_t = s_i) = \max_{s_1 \leq z_{t-1} \leq s_N} (\delta_{t-1}(z_{t-1}) \cdot a_{z_{t-1}, z_t}) b_{z_t}(x_t)$

$\psi_t(z_t) = \arg \max_{s_1 \leq z_t \leq s_N} \delta_{t-1}(z_{t-1}) \cdot a_{z_{t-1}, z_t}$

iii. Max over all end points.

$p^* = \max_{s_1 \leq z_T \leq s_N} \delta_T(z_T) \rightarrow$ final most likely path that achieves max end state prob.

$s_1 \leq z_T \leq s_N$

$z_T^* = \arg \max_{s_1 \leq z_T \leq s_N} \delta_T(z_T)$

$s_1 \leq z_T \leq s_N$

most likely end state.

iv. Path reconstruction

$z_{t+1}^* = \psi_{t+1}(z_{t+1}^*)$

- $O(N^2 t)$: computational complexity

- Viterbi is also used in decoders

$\alpha_t(s_i) \rightarrow$ sum of path prob. of all paths fronting 1 to time t that are at time t in state s_i

$\beta_t(s_i) \rightarrow$ sum of prob. of all paths from time t to time T that are at time t in state s_i

$\psi_t(z_i) \rightarrow$ caches predecessor state to reconstruct path.

$\delta_t(s_i) \rightarrow$ prob. of single most likely path from time t to time T that is at time t in state s_i

Algo 3: Training

Baum-Welch formulae = EM algorithm

Expectations

- responsibilities are estimations for current state/state transition & associated expectations
- idea: FWD A BWD algo marginalise most parameters to isolate param. of interest.

Maximization:

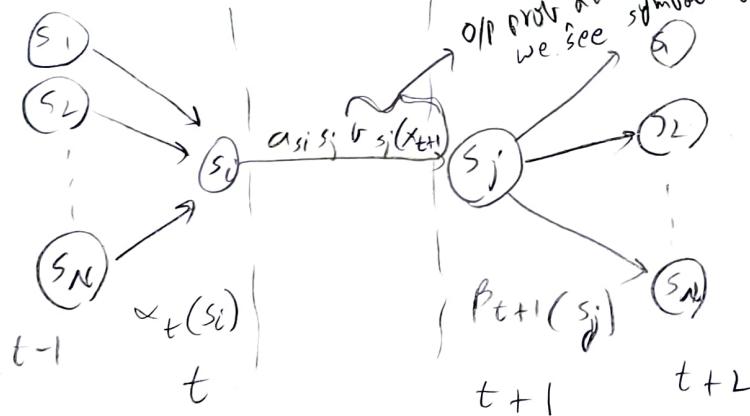
updates model param.

idea: count relative freq. converges to local opt.

x_t is a position in an i/p seq. x_1, \dots, x_T

However training is done over a dataset/ many i/p sequences.

- state seq. is hidden
- for training we have to update individual state transitions α_{ij}, β_j
- done by accessing $\alpha_t(z_t = s_i) + \beta_{t+1}(z_{t+1} = s_j)$



$$\text{ex: } \alpha_{15,7} \quad \alpha_t(15) \quad \beta_{t+1}(7)$$

STEP E likelihood of state transition $s_i \rightarrow s_j$ at time t

$$\begin{aligned} \xi_t(i, j) &= P(z_t = s_i, z_{t+1} = s_j) \\ &= \frac{\alpha_t(s_i) \alpha_{is_j} \beta_{sj}(x_{t+1}) \beta_{t+1}(s_j)}{\sum_{z_t=s_i} \sum_{z_{t+1}=s_j} \alpha_t(z_t) \cdot \alpha_{z_t z_{t+1}} \beta_{z_{t+1}}(x_{t+1}) \cdot \beta_{t+1}(z_{t+1})} \end{aligned}$$

likelihood for being in s_j at time t : marginalize s_j

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j)$$

expected no. of transitions $s_i \rightarrow s_j$ at all times: $\sum_{t=1}^T \xi_t(i, j)$
" " " " times s_i is visited: $\sum_{t=1}^T \gamma_t(i)$

$\hat{\pi}_{s_i} = \gamma_1(i) \xrightarrow{\text{expected}} \# \text{times } s_i \text{ in time 1 over all training data}$
update $\bar{\alpha}_{s_i s_j}$ (state transition prob.)

$$\bar{\alpha}_{s_i s_j} = \frac{\sum_{t=1}^T \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)} = \frac{\text{exp. # of trans. } s_i \rightarrow s_j}{\text{exp. # of times in state } s_i}$$

update O/P prob.

$$\bar{\nu}_{s_i(Ov)} = \frac{\sum_{t=1}^T \gamma_t(i) \mathbb{I}(x_t, ov)}{\sum_{t=1}^T \gamma_t(s_i)} = \frac{\text{exp. # times in } s_i \text{ obs. or}}{\text{expected # times in } s_i}$$

HMM = ergodic if every state can be reached from every other state

[speech is not ergodic]

Left \rightarrow Right HMMs are mostly used. [$\alpha_{s_i s_j} = 0$ if $i > j$]
(only Fwd & self loops allowed)

$$A = \begin{bmatrix} 0.1 & 0.6 & - & - \\ 0 & - & - & - \\ 0 & - & - & - \\ 0 & - & - & - \end{bmatrix}$$

- GMM & HMM training work well together
 - both probabilistic models
 - " trained using EM.
- [can nest these eq. together]
- good init. can improve results.

Markov Random Fields

- Undirected graphical model
- used for labeling tasks: segmentation, stereo matching, multi-view stitching.
- idea: consider unknown labels as a field of RV with conditional independance,

Joint dist. & objective

$$p(\underbrace{x_1 \dots x_N}_{\text{observations}}, \underbrace{z_1 \dots z_N}_{\text{hidden variables}}) = p(x_1 \dots x_N | z_1 \dots z_N) p(z_1 \dots z_N)$$

MRF inference seeks posterior of hidden var.

$$p(z_1 \dots z_N | x_1 \dots x_N) = \frac{p(x_1 \dots x_N | z_1 \dots z_N) p(z_1 \dots z_N)}{p(x_1 \dots x_N)}$$

$$= \frac{p(x_1 \dots x_N | z_1 \dots z_N) p(z_1 \dots z_N)}{\sum_{z_1} \sum_{z_N} p(x_1 \dots x_N | z_1 \dots z_N) p(z_1 \dots z_N)}$$

marginalisation
of numerators

here evidence needs exponential calculations as we have to sum over all possible combinations of variables, hence we will seek Max. likelihood sol.

then - we will show the denominator is const. WRT z .

- Likelihood & prior must also be simplified.

- Ex: in images $x_i \rightarrow i/p$ pixel

$z_i \rightarrow$ label we want to assign

for low res. images 256×256 pixels $\rightarrow 2 \cdot 2^{16}$ variables.

Joint distr. is too complex [look at graph in slide]

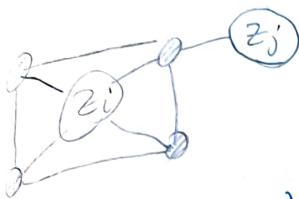
How A MRF make model tractable with simple dependency

Factorization of joint dist

1. Each obs. depends only on 1 hidden var.
2. Each z depends only on its neighbours i.e. $z_i \perp\!\!\!\perp z_j | N(z_i)$
for all $z_j \notin N(z_i)$. [$N(z_i) \rightarrow$ set of neighbours of z_i]
 z dependent on neighbours but independent of non-neighbours.

(1) neighbourhood of z_i should be split in pairs for tractability such that prior $p(z_1, \dots, z_N)$ decomposes to

$$p(z_1, \dots, z_N) = \prod_{z_i \in N(z_j)} p(z_i, z_j) \prod_{z_i \in N(z_j)} p(z_i, z_1)$$



point 2a → highlights neighborhoods determine factorization.

Undirected graphs factorize into maximal cliques.

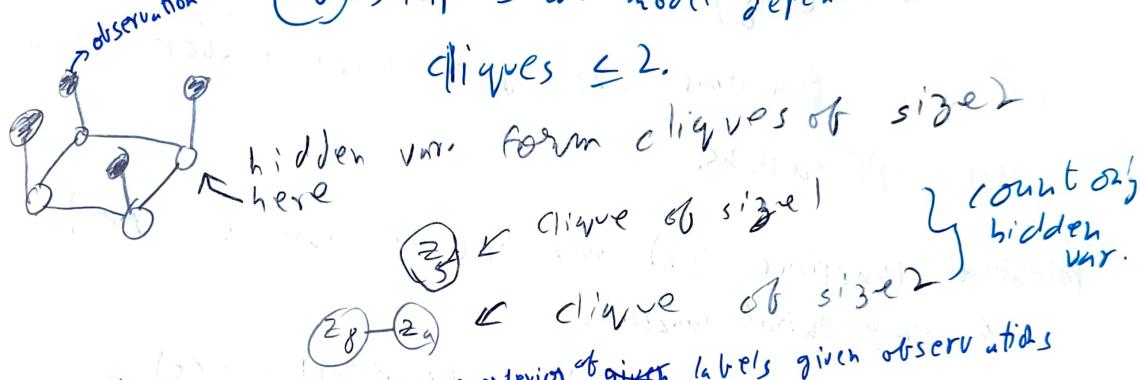
Clique → fully connected subgraph.

dependency structure.

$$p(a, b, c, d, e, f, g) = p(a, b, c) p(c, d) p(b, f) p(d, e, f, g)$$

maximal fully
connected subgraphs

(2b) → implies we model dependencies with cliques ≤ 2 .



These assumptions result in posterior of given labels given observations

$$p(z_1, \dots, z_N | x_1, \dots, x_N) = \prod_{i=1}^N p(x_i | z_i) \prod_{z_i \in N(z_j)} p(z_i, z_j)$$

$$\sum_{z_1, \dots, z_N} \prod_{i=1}^N p(x_i | z_i) \prod_{z_i \in N(z_j)} p(z_i, z_j)$$

Inference for MRF is to find [Max Likelihood] inference : tractable

$$z^* = \arg \max_z p(z_1, \dots, z_N | x_1, \dots, x_N) \quad - \text{learning: intractable}$$

↳ to find $p(x_i | z_i)$
+ $p(z_i, z_j)$
these are difficult to handcraft

with learning being intractable & difficult to handcraft inference is impossible.

Trick: use potentials instead of distributions.

Hammersley-Clifford Theorem

- define $p(x_i|z_i) \wedge p(z_i, z_j)$ optimally by using energy terms
- Hammersley-Clifford theorem allows to move from MRF to Gibbs Random Field (GRF) which wraps energy terms in exponentials.

$\psi(c) \rightarrow$ potential functions over cliques of var. C.

$$\psi(c) \geq 0 \quad [\text{non-ve}]$$

MRF = GRF [consisting of $\psi(c)$]

Instead of specifying $p(x_i|z_j) \wedge p(z_i, z_j)$ set $\psi(c)$ where $C = \{z_i\}$ are maximal cliques of hidden variables.

Potential functions are often more easier to choose, in many applications.

Potential functions $\psi(c)$ over set of nodes $C = \{z_i\}$ are set of exponential functions

$$\psi(c) = \exp(-E(c)); \quad E(c) \rightarrow \text{energy function}$$

lower energy = closer to goal

\therefore factorized MRF is

$$\prod_{i=1}^N p(z_1, \dots, z_N | x_1, \dots, x_N) = \frac{1}{Z} \prod_c \psi(c)$$

where $Z = \sum_c \prod_c \psi(c) \rightarrow$ called partition function

$Z \rightarrow$ enables use of any non-negative func. for potentials

$Z \rightarrow$ intractable, but is const - when minimizing $\prod c$
 \therefore it can be ignored.

Inference with potentials

$$z^* = \arg \max_z p(z_1, \dots, z_N | x_1, \dots, x_N) = \arg \max_z \frac{1}{Z} \prod_c \psi(c)$$

Inference via Gibbs Sampling (MC MC based)

1. select node z_i (randomly/follow a pattern)
2. get current labels from all neighbours (Markov Blanket)
3. Assign a new label to z_i with minimal energy.
4. goto 1 for set iterations/until convergence.

→ go for min. energy / a/m for most likely sol. instead of sampling any sol.

This makes sol. locally optimal.

Graph cuts find globally optimal sol.

Unary potentials: potentials with energy terms $E(x_i, z_i)$

Pairwise " : " " " " $E(z_i, z_j)$

Optimization analogy: unary potentials the data term & the pairwise potentials are regularizer.

Example ~~Belief propagation~~ Denoising binary img

every func. for
obs-hidden connections $\exists z_i$

$$E(x_i, z_i) = -\eta x_i z_i ; \quad \begin{aligned} \eta &\rightarrow \text{weight constant} \\ x_i &\rightarrow i^{\text{th}} \text{ ip pixel} \\ z_i &\rightarrow i^{\text{th}} \text{ hidden variable} \end{aligned}$$

energy func. for hidden-hidden connections

$$E(z_i, z_j) = -\beta z_i z_j \quad \beta \rightarrow \text{weight const.}$$

-Value range of labels + ips constraints the potential

returned pixel label
is the denoised
image

if it's a binary image then

there are 2 possibilities

$$\text{identical labels} \rightarrow E(x_i = -1, z_i = -1) = -\eta$$

so we want to encourage
is lowest energy hence
this is favoured

$$\text{diff. labels} \rightarrow E(x_i = -1, z_i = 1) = +\eta$$

more

same logic holds for hidden-hidden connections

data energy term

$$E(z_i) = h \sum_{z_i} + \beta \sum_{z_i, z_j \in N(z_i)} z_i z_j - \gamma \sum_{(x_i, z_i)}$$

first term to
decide a label

non-binary labeling

Unary term could be $\gamma \|x_i - y_i\|_2^2$. low energy sol. \hat{x}
that are close to y_i .

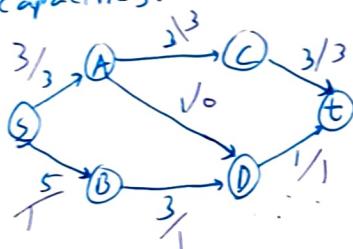
pairwise energy term could be $\beta \|z_i - z_j\|_2^2$. low
energy sol. where neighbouring pixels y_i, y_j are more
similar.

Inference via Min-cut & Max flow.

What is max-flow?

- combinatorial std. problem
- solved in polynomial time.

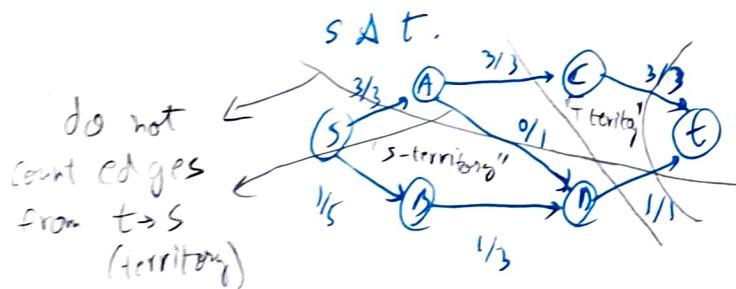
Determine max throughput of water from source(s) to sink(t)
are tube capacities.



→ pumped units
Max flow = 4

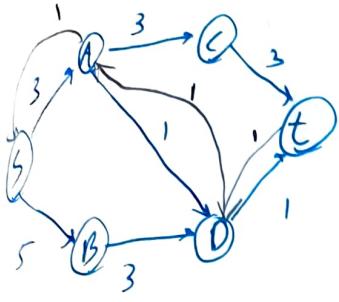
What is min-cut?

- seeks smallest sum of edges to disconnect s & t.
- select minimal set of naked out edges that disconnects



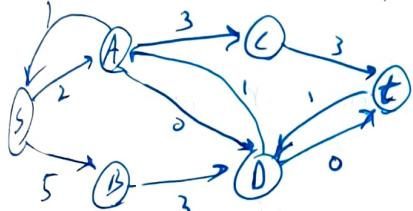
Ford fulkerson Algo:

finds out Max flow

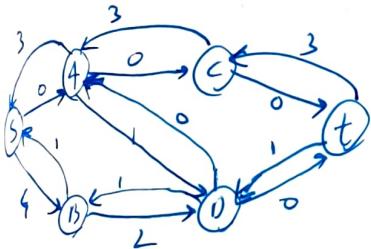


shortest path: $S \rightarrow A - D - T$ [S]
path flow = 1

capacity of 3 becomes 2
() " 0



shortest path $S \rightarrow A \rightarrow C \rightarrow T$
path flow = 2



MRF inference via Graph cuts

MRF \rightarrow GRF \rightarrow Energy minimisation

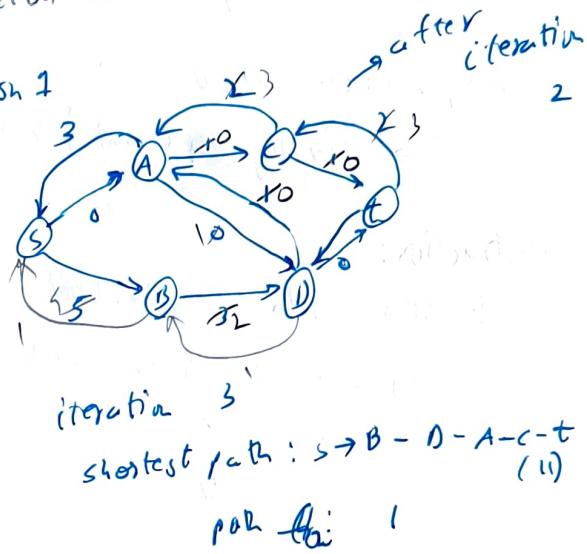
Objective of MRF inference task: Find optimal label assignments

Idea of MRF inference via graph cuts: transform maximum likelihood task to equivalent min cut task.

Procedure: construct directed graph with edge weights such that graph's min cut provides sol. to max likelihood task.

1. greedily search shortest path
2. max-out flow capacity on that path, reduce edge weights
3. Introduce backward edges to undo greedy edges to undo greedy selection from greedy selection

after iteration 1



iteration 3

shortest path: $S \rightarrow B - D - A - C - T$ (1)
path flow = 1

total flow through iterations = $1 + 2 + 1 = 4$

min cut = $D \rightarrow T$

either $(S-A, A-C, C-T)$

Requirements for graph construction / constraints for algo

1. Binary labels

2. Max. clique size of 2

3. Pairwise energy terms satisfy submodularity condition

$$E(0,0) + E(1,1) \leq E(0,1) + E(1,0)$$

Under these 3 constraints algo. finds

- globally optimal sol.

- in polynomial time.

& expansion used to extend this method to non-binary labeling
↳ locally optimal [but with guaranteed margin around global opt]

Construction:

1. start with neighborhood relationship of hidden var z_i .

2. Encode optimization prob.

a. add source node s , identified with label 0

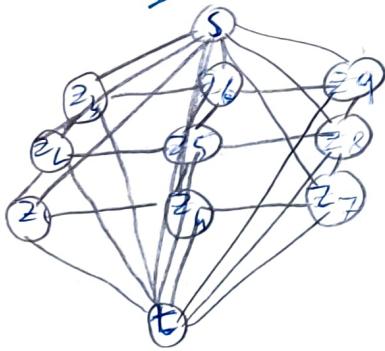
b. " sink " t , " " " " 1

c. connect nodes with appropriately chosen weight

3. calc min. cut b/w s & t

4. Assign "0" to nodes connected to s

"1" to nodes connected to t



Choosing edge weights: Additivity of graphs

Equivalence b/w MRF inference & graph cut ensured by construction of edge weights.

Every energy term can be transformed to graph snippet because energy minimization via min-cut is **homomorphic**.

i.e. if g_1 encodes $\min(E_1)$ & g_2 encodes $\min(E_2)$
then $(g_1 \cup g_2)$ encodes $\min(E_1 + E_2)$

if we encode each single energy term to a graph such that graph's min cut encodes the min energy, then also min cut of sum/union of all individual graphs will provide min energy of sum of terms.

Sol: to min energy of sum of terms.

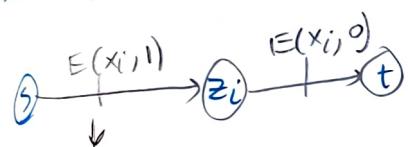
if $E(x_i, 1) > E(x_i, 0)$



if $E(x_i, 0) > E(x_i, 1)$



Encoding unary term



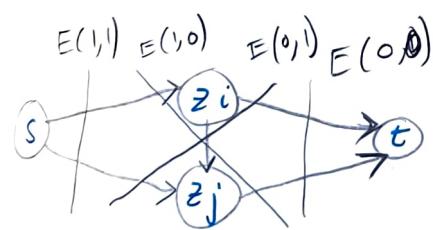
smaller graph
 $s=0 \wedge t=1$

cut b/w s & z_i assigns label $z_i = t = 1$

cost is $E(x_i, 1) \rightarrow$ relates observation x_i to 1

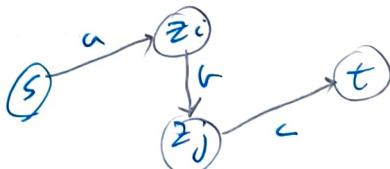
smaller graphs
computationally cheaper
but shift minimum's offset.

Encoding pairwise term



$s=0 \wedge t=1$
smaller graph

if $E(1,0) > E(0,0)$
 $E(1,0) > E(1,1)$



a, b, c obtained from

$$a+b = E(1,1)$$

$$c+b = E(0,0)$$

$$(r+b = E(0,1))$$

$$a+c+b = E(1,0)$$

$b \rightarrow \text{const}$ offset.

α -expansion

1. select 1 specific label α from the label set
2. fix hidden variables that already have label α
3. seek lowest energy labeling that switch atleast 1 other hidden var. to α
4. keep expanded α labeling if its energy is lower than current energy.