## Search Problems Definition

state space $S$, initial state $s_0 \in S$, actions $A(s)$ of state $s$, transition model/function Result$(s,a)$, goal test $G(s)$, action cost $c(s,a,s')$

**State space graph:** Node are states, Arcs represent transitions, goal test is a set of goal nodes

**Search Tree:** root node is start state, children correspond to successors, nodes show states, correspond to plans that achieve that states. Node in Search Tree = Path in the state space

## Uninformed Search:

| Criterion | Strategy | Frontier Implementation | Action Costs | Complete | Optimal | Time | Space | Comment |
|---|---|---|---|---|---|---|---|---|
| Depth-First | expand deepest node first | LIFO stack | $=c$ | No | No | $O(b^m)$ | $O(bm)$ | $m$ is tier num |
| Breadth-First | expand shallowest node first | FIFO queue | $=c$ | Yes | Yes | $O(b^d)$ | $O(b^d)$ | |
| Iterative Deepening | DFS with increasing depth limit | | $=c$ | Yes | Yes | $O(b^d)$ | $O(bd)$ | |
| Uniform-Cost | expand lowest $g(n)$ $g(n)$:root → n cost | priority queue sorted by $g(n)$ | $\geq \epsilon$ | Yes | Yes | $O(b^{C^*/\epsilon})$ | $O(b^{C^*/\epsilon})$ | Solution cost $C^*$, minimal arc cost $\epsilon$ |

**Heuristic:** A function that estimates how close a state is to a goal, $h(n)$

**Greedy Search:** Strategy: expand a node that you think is closest to a goal state. Not Complete. Not optimal. Worst case: badly guided DFS.

**A\* search:** Strategy: UCS + Greedy Search, sort by $f(n) = g(n) + h(n)$.

**Tree search optimal & Admissible Heuristics:** $0 \leq h(n) \leq h^*(n)$ estimate cost < real cost → Provide optimal property for tree search

**Graph search optimal & Consistency of Heuristics:** $h(n) - h(n') \leq c(n,a,n')$ estimate cost < real cost for each arc [Consistency → admissible] → Provide optimal property for graph search.

**Creating Heuristics:** Using the answer of a relaxed (constraint) problem; mixing heuristics

**Constraint Satisfaction Problem(CSP)** $P = (X,D,C)$, Variables: $X = \{X_1,...,X_n\}$, Domains: $D = \{D_1,...,D_n\}$, each domain $D_i = \{v_1,...,v_k\}$ for variable $X_i$; Constraints $C$, allowable combinations of values(binary CSP: only binary constraint); Assignment: $\{X_i = v_i, X_j = v_j\}$ (complete: every variable is assigned); Solution: a consistent and complete assignment

**Solution: 1.** Backtracking Search = DFS + variable-ordering + fail-on-violation(With Dynamic Ordering: (a)MCV Choose variable that has the fewest consistent values. (b)LCV Order values of selected $X_i$ by decreasing number of consistent values of neighboring variables.)

**Structured CSP problems:** No loop constraint graph can be solved in linear time: $O(nd^2)$

**Arc Consistency:** Consistent arc $X_i \rightarrow X_j$: for every $x_i \in \text{Domains}_i$, exists $x_j \in \text{Domain}_j$ that is consistent; Enforce arc consistency: Reduce $\text{Domain}_i$

**AC-3 Algorithm:** if $X_i$ lose a value, the arc points to $X_i$ need to be rechecked. **Implementation:** Use a queue to store arcs to be checked. Time complexity: $O(c \cdot d \cdot d^2)$, each arc could be insert into queue for most $d$ times, every time it takes $d^2$ time [Note: it isn't always effective]

Called when the domain of $X_t$ is reduced.

```
function AC-3(X_t, Domains) returns false if an inconsistency is found and true otherwise
  initialize queue with all arcs (X_k, X_t) for X_k in Neighbours(X_t)
  while queue is not empty do
    (X_t, X_j) ← RemoveFirst(queue)
    if EnforceArcConsistency(Domains, X_t, X_j) then  The domain of X_t is reduced.
      if size of Domains_t = 0 then return false
      for each X_k in Neighbours(X_t) do
        add (X_k, X_t) to queue
  return true
```
Constraint propagation: $X_k \rightarrow X_t \rightarrow X_j$

## Local search

**Local search:** faster, memory efficient, but incomplete and suboptimal. General Idea: improves a single option until you can't make it better.

**Hill Climbing:** start wherever, REPEAT: move to the best neighbor. Problems: local maxima(stuck), plateaus(lost) Variants: Stochastic(randomly choose move) Random-restart(randomly restart)

**Simulated Annealing:** Gradually reduce temperature; higher temperature, more bad moves allowed, shake the system out of its local best(bad move accept rate:$\exp \Delta E/T$)

**Local Beam Search:** greedily keep $k$ (best) states. No optimal property. Running time $O(nkb \log kb)$ Variant: Stochastic (choosing successors randomly, avoid lack of diversity.)

## Standard Game Formulation

Initial state $s_0$. Players Player$(s)$ who is playing $s$. Actions: Actions$(s)$. Transition model: Result$(s,a)$. Terminal test: Termnal_Test$(s)$. Utility function: Utility$(s,p)$.

**Minimax Algorithm**: a state-space tree search(DFS), choose the action leading to state with best achievable utility against an optimal adversary.

**MAX nodes**: Under our control
$V(s) = \max_{a \in \text{Actions}(s)} V(\text{Result}(s,a))$

**MIN nodes**: Under opponent's control
$V(s) = \min_{a \in \text{Actions}(s)} V(\text{Result}(s,a))$

**Generalized Minimax:** Each player maximizes its own component

### Alpha-Beta Pruning

```
function MAX-VALUE(s, α, β) returns a v
  if Terminal-Test(s) then return Utility(s)
  initialize v = -∞
  for each successor s' of state s:
    v = max(v, MIN-VALUE(s', α, β))
    if v ≥ β
      return v
    α = max(α, v)
  return v
```

**MCTS(Monte Carlo Tree Search):** use rollout

**Naïve/Pure:** do $N$ complete random rollouts, choose the one with best percentage win

**Better Rollout Policy:** Exploitation(focus on more promising moves) & Exploration (focus on more uncertain nodes)

**UCB I Formula(Upper confidence bound):**

value estimate $v_i + C \times \sqrt{\dfrac{\ln(N)}{n_i}}$

tunable parameter; total number of trials; num trials for arm $i$

Exploit: prefers higher payoff arm
Explore: prefers less played arm

## MCTS Steps:

**1. Selection**, used for nodes we have seen before, pick according to UCB

**2. Expansion**, used when we reach the frontier, add one node per playout

**3. Simulation**, used beyond the search frontier; no UCB, just play randomly

**4. Backpropagation**, after reaching a terminal node, update value and visits for states expanded in selection and expansion

• Pros:
- Grows tree asymmetrically, balancing expansion and exploration
- Unaffected by branching factor
- Easy to adapt to new games
- Heuristics not required, but can also be integrated
- Anytime algorithm: can finish on demand
- Trivially parallelizable

• Cons:
- Can't handle extreme tree depth
- Requires ease of simulation, massive computation resources
- Relies on random play being "weakly correlated"
- Many variants, need expertise to tune
- Theoretical properties not yet understood

**Components of Learning**: Input: $x \in \mathbb{R}^d$, Output: $y = \{0,1\}$, Target Function: $f: \mathcal{X} \rightarrow \mathcal{Y}$, Data: $(x_1,y_1),...,(x_n,y_n)$, [Supervised Learning], Hypothesis: $h: \mathcal{X} \rightarrow \mathcal{Y}$, Space: $\mathcal{H} = \{h\}$, Learning algorithm $\mathcal{A}$ to find the best $h \approx f$ Loss function: $l: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$, measure the difference

| | Supervised Learning | Unsupervised Learning |
|---|---|---|
| Discrete Output | Classification 分类 | Clustering 聚类 |
| Continuous Output | Regression 回归 | Embedding 嵌入 |

| | | Ground Truth | |
|---|---|---|---|
| | | Actual Positive | Actual Negative |
| Predicted | Predicted Positive | TP True Positive | FP False Positive |
| | Predicted Negative | FN False Negative | TN True Negative |

**Precision** = $\frac{TP}{TP+FP}$ : How many selected items are relevant

**Recall** = $\frac{TP}{TP+FN}$ : How many relevant items are selected

$TPR = \frac{TP}{TP+FN}$, $FPR = \frac{FP}{FP+TN}$

Sensitivity = Recall = TPR = $P(\hat{Y}=1|Y=1)$
Specificity = 1 − FPR = $P(\hat{Y}=0|Y=0)$

**KNN**: **When to use:** 1. Few attributes per instance (expensive computation); 2. Lots of training data (curse of dimensionality)

**Advantages:** 1. Agnostically learn complex target functions; 2. Do not lose information (store original data);3. Data number can be very large (big pro!); 4: Class number can be very large (biggest pro!)[All other ML algorithms may fail here!]

**Disadvantages**: 1. Slow at inference time (acceleration a must) 2. Fooled easily by irrelevant attributes (feature engineering crucial)

**Z-score normalization:**
For each feature dimension $j$, compute based on its samples: 1.

**Mean** $\mu_j = \frac{1}{N}\sum_{i=1}^N x_{ij}$, **Variance** $\sigma_j = \sqrt{\frac{1}{N}\sum_{i=1}^N (x_{ij} - \mu_j)^2}$, 2:

**normalize** the feature into a new one $\hat{x}_{ij} \leftarrow \frac{x_{ij} - \mu_j}{\sigma_j}$

## SVC

• Apply Representer theorem, we can kernelize Soft-SVM as follows:
$$\min_{\alpha} \frac{1}{2}\alpha^T K\alpha + \frac{C}{n}\sum_{i=1}^n \max\left\{0, 1 - y_i \sum_{j=1}^n \alpha_j k(x_i, x_j)\right\}$$

- $\alpha_j$ is the weight of each reference point $x_j$ to the prediction of $x_i$
- It is actually a Primal Form with kernel functions.

**Hard-margin Support Vector Machine (SVM)**
$$\min_{w,b} \frac{1}{2}\|w\|_2^2$$
$$\text{s.t. } y_i(w \cdot x_i + b) \geq 1, 1 \leq i \leq n$$

**Soft-margin Support Vector Machine (SVM)**
$$\min_{w,b,\xi} \frac{1}{2}\|w\|_2^2 + \frac{C}{n}\sum_{i=1}^n \xi_i$$
$$\text{s.t. } y_i(w \cdot x_i + b) \geq 1 - \xi_i$$
$$\xi_i \geq 0, 1 \leq i \leq n$$

• Define Hinge loss $\ell(f(x),y) = \max\{0, 1 - yf(x)\}$
- For the linear hypothesis: $\ell(f(x),y) = \max\{0, 1 - y(w \cdot x + b)\}$
• Theorem: Soft-SVM is equivalent to a Regularized Risk Minimization:
$$\min_{w,b} \frac{1}{2}\|w\|_2^2 + \frac{C}{n}\sum_{i=1}^n \max\{0, 1 - y_i(w \cdot x_i + b)\}$$

## Stochastic Gradient Descent

• In each iteration $t$ $(\leq T)$:
- Randomly sample a minibatch of $m \ll n$ points $\{(x_i, y_i)\}_{i=1}^m$
- Set $J^t(\theta^t) = \frac{1}{m}\sum_{i=1}^m \ell(\theta^t; x_i, y_i)$
- Compute gradient on minibatch: $\Delta^t = \nabla_\theta J^t(\theta^t)$
- Update parameters with learning rate $\eta$: $\theta^{t+1} = \theta^t - \eta \Delta^t$

## Logistic Regression

2-class Classification:

Hypothesis Space $h_w(x) = w \cdot \Phi(x)$

$\sigma(w \cdot z)$ + discretization

{0,1}
Discrete Label Space $\mathcal{Y} = \mathbb{R}$
Output Space $\mathcal{Y} = \mathbb{R}$

Feature Engineering; Polynomial Basis Function $\Phi$

Input Space $\mathcal{X} = \mathbb{R}^d$; High-dim Space $\mathcal{Z} = \mathbb{R}^s$

$$\hat{e}(w) = -\sum_{i=1}^n \{y_i \log \sigma(h_w(x)) + (1-y_i)\log[1 - \sigma(h_w(x))]\} + \lambda\Omega(w)$$

### Linear Regression

Hypothesis Space $h_w(x) = w \cdot \Phi(x)$

Gradient Descent
$\nabla_w \hat{e}(w) = 2Z^T(Zw - y) + 2\lambda w$

Loss Function
$$\min_w \sum_{i=1}^n (h_w(x_i) - y_i)^2 + \lambda\Omega(w)$$

Analytic Solution
$w = (Z^T Z + \lambda I)^{-1} Z^T y$

Well Done! A simple but complete model!

• Softmax function normalizes multiple outputs in a probability vector:

### Softmax
$$p(y=i|x) = \frac{\exp(w_i^T x)}{\sum_{r=1}^C \exp(w_r^T x)}$$
Sum over all classes.

**Hypothesis of decision tree:** 1. Decision tree divide the feature space into axis-parallel rectangles 2. each rectangular region is labeled with a specific label

**Entropy(Use in ID3 and C4.5):** $H(\mathcal{D}) = -\sum_{k=1}^k \frac{|\mathcal{C}_k|}{|\mathcal{D}|} \log \frac{|\mathcal{C}_k|}{|\mathcal{D}|}$ k

Maximizing **Information Gain:** $H(\mathcal{D}_1 \cup \mathcal{D}_2) - \frac{|\mathcal{D}_1|}{|\mathcal{D}|}H(\mathcal{D}_1) - \frac{|\mathcal{D}_2|}{|\mathcal{D}|}H(\mathcal{D}_2)$

### ID3-Algorithm

数据集 Class label 特征集
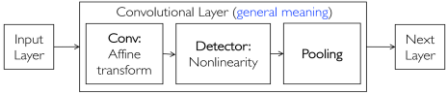ID3(Examples, Target_attribute, Attributes)

Stop Criteria
• create a Root node for the tree; assign all Examples to Root
• if all Examples are positive, return the single-node tree Root, with label=+;
• if all Examples are negative, return the single-node tree Root, with label=−;
• if Attributes is empty, return the single-node tree Root, with label = the most common value of Target_attribute in Examples
• otherwise // Main loop:
  $A \leftarrow$ the attribute from Attributes that best* classifies Examples; 选出IG最大的
  the decision attribute for Root ← $A$; $O(dn)$ in each layer
  for each possible value $v_i$ of $A$:
    add a new tree branch below Root, corresponding to the test $A = v_i$;
    let Examples$_{v_i}$ be the subset of Examples that have value $v_i$ for $A$;
    if Examples$_{v_i}$ is empty: 样本没有, 测试可能有的
      below this new branch add a leaf node with label = the most common value of Target_attribute in Examples;
    else
      below this new branch add the subtree depth min$(d, \log n)$
      ID3(Examples$_{v_i}$, Target_attribute, Attributes\{$A$});
• return Root;
  * The best attribute is the one with the highest information gain.
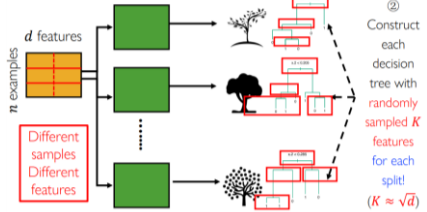
**CNN: 1.** Local Assumption: local information is enough for recognition ➔ only connect local neurons **2.** Shift Invariance Assumption: if a feature is useful at spatial position ➔ share the weight of sliding windows, **Size** $\text{height}_{new} = [(\text{height} - \text{filter} + 2 \cdot \text{pad})/\text{stride}] + 1$ $\text{weight}_{new} = [(\text{width} - \text{filter} + 2 \cdot \text{pad})/\text{stride}] + 1$

Convolutional Layer (general meaning)
Input Layer → Conv: Affine transform → Detector: Nonlinearity → Pooling → Next Layer

**RNN: Local Dependency Assumption:** The sequential information of all previous timestamps can be encoded into one hidden representation. **Temporal Stationarity Assumption:** If a feature is useful at time $t_1$, then it should also be useful for all time stamps $t_2$
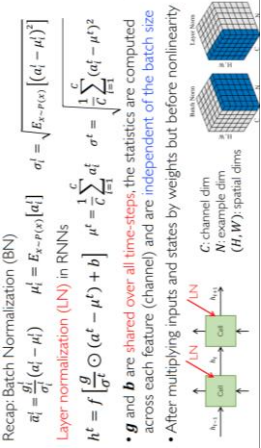
### Random Forest

• Bootstrap (自助法): randomly draw datasets with replacement from the training data, with the same sample size as the original training set
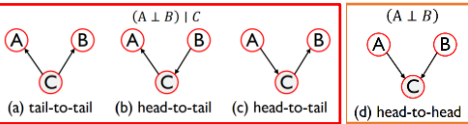① Bootstrap samples from training data
② Construct each decision tree with randomly sampled $K$ features for each split!
$(K \approx \sqrt{d})$

Different samples Different features

• **Probabilistic Reasoning = Modeling + Inference**
- $X = \{x_1, ..., x_D\}$ is a set of $D$ random variables.
- Query set $R$ and condition set $C$ are subsets of $X$.
• Modeling: How to specify a joint distribution $p(x_1, ..., x_D)$ compactly?
• Inference: How to compute $p(R \mid C)$ efficiently?
• A Bayesian network is a directed acyclic graph (DAG) that specifies a joint distribution as a product of local conditional distributions, one for each node: $p(x_1, ..., x_K) = \prod_{s=1}^{K} p(x_s | x_{\Gamma(s)})$ where $\Gamma(s)$ denotes the set of parents of $x_s$.
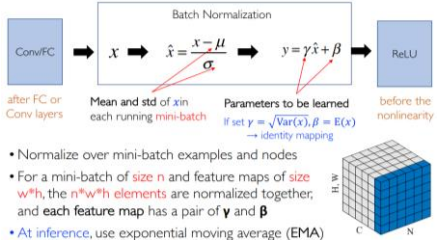
### EM for GMM

E-step. Evaluate all responsibilities using current parameters:
$$\gamma_i^j = \frac{\pi_j \mathcal{N}(x_i|\mu_j, \Sigma_j)}{\sum_{c=1}^{k} \pi_c \mathcal{N}(x_i|\mu_c, \Sigma_c)}$$
M-step. Re-estimate the parameters using the responsibilities:
$$\pi_c^{new} = \frac{n_c}{n}$$
$$\mu_c^{new} = \frac{1}{n_c} \sum_{i=1}^{n} \gamma_i^c x_i$$
$$\Sigma_c^{new} = \frac{1}{n_c} \sum_{i=1}^{n} \gamma_i^c (x_i - \mu_c^{new})(x_i - \mu_c^{new})^T$$
$$n_c = \sum_{i=1}^{n} \gamma_i^c$$

### Batch & Layer

Recap: Batch Normalization (BN)
Layer normalization (LN) in RNNs
• $g$ and $b$ are shared over all time-steps, the statistics are computed across each feature (channel) and are independent of the batch size
• After multiplying inputs and states by weights but before nonlinearity

### Conditional Independence / Marginal Independence

$(A \perp B) \mid C$
(a) tail-to-tail (b) head-to-tail (c) head-to-head
$(A \perp B)$
(d) head-to-head

判断条件独立: **1.** 找到 C 及 C 的所有祖先，断开**不是** h-to-h 的边。**2.** 对其他的点，断开**是** h-to-h 的边。
结果：如果不联通，则条件独立。

Batch Normalization
Conv/FC → $x$ → $\hat{x} = \frac{x - \mu}{\sigma_c}$ → $y = \gamma \hat{x} + \beta$ → ReLU
after FC or Conv layers
Mean and std of $x$ in each running mini-batch
Parameters to be learned If set $\gamma = \sqrt{Var(x)}, \beta = E(x)$ → identity mapping
before the nonlinearity

• Normalize over mini-batch examples and nodes
• For a mini-batch of size n and feature maps of size w*h, the n*w*h elements are normalized together, and each feature map has a pair of $\gamma$ and $\beta$
• At inference, use exponential moving average (EMA)

---

• How to infer a node with variable $x_i$ from the remaining variables $x_{i \neq i}$ in directed models ? $p(x_i \mid x_{\{j \neq i\}})$
• Find $S_1 \in x_{\{j \neq i\}}$ which contents **Markov blanket**
• It means that $S_1$ contains all information one needs to infer $x_i$.
$$p(x_i \mid x_{\{j \neq i\}}) = p(x_i \mid S_1 \sqcup (x_{\{j \neq i\}} \setminus S_1)) = p(x_i \mid S_1)$$
• $S_1$ is called the Markov blanket of $x_i$. $\quad x_i \perp (x_{\{j \neq i\}} \setminus S_1) \mid S_1$
• A Markov boundary is the minimal Markov blanket. In a Bayesian network, it includes parents, children and the other parents of all of its children (co-parents). $p(x_i \mid x_{\{j \neq i\}}) = p(x_i \mid x_{pa} \cup x_{ch} \cup x_{co-pa})$
• Marginal inference · What is the probability of a given variable in our model after we sum everything else out?
$$p(y = 1) = \sum_{x_1} \sum_{x_2} ... \sum_{x_n} p(y = 1, x_1, x_2, ..., x_n)$$

• Algorithm (Variable Elimination, VE)
• For each variable $X_i$ (ordered according to $O$):
1. Multiply all factors $\phi_i$ containing $X_i$ ; 乘法
2. Marginalize out $X_i$ to obtain a new factor $\tau$ ; 加法
3. Replace the factors $\phi_i$ with $\tau$. 换掉
• Consider a general distribution $P(X, Y, E)$ over sets of:
- Query variables $Y$; - Observed evidence variables $E$;
- Unobserved variables $X$.
$$P(Y \mid E = e) = \frac{P(Y, E = e)}{P(E = e)} \quad \text{Apply VE Algorithm!}$$
• We can select the elimination ordering $X \to Y \to E$ to obtain $P(Y, E = e)$ and $P(E = e)$ in a single turn of VE algorithm.

**Learn:** Give data $\mathcal{D}$, find "best" parameter $\theta$
**Inference:** Give data $x$ and parameter $\theta$, find latent variables $z$'s distribution.

**Parametric model** $\{p(z; \theta) \mid \theta \in \Theta\}$
A sample dataset $\mathcal{D} = (z_1, ..., z_N)$
Likelihood of $\hat{\theta} \in \Theta$ for sample $\mathcal{D}$ is:
$$p(\mathcal{D}; \hat{\theta}) = \prod_{n=1}^{N} p(z_n; \hat{\theta})$$
Log likelihood due to numerical reason:
$$\log p(\mathcal{D}; \hat{\theta}) = \sum_{n=1}^{N} \log p(z_n; \hat{\theta})$$
**MLE(Maximum Likelihood estimator):**
$$\hat{\theta} \in \arg\max_{\theta \in \Theta} \log p(\mathcal{D}; \theta)$$
$$= \arg\max_{\theta \in \Theta} \sum_{n=1}^{N} \log p(z_n; \theta)$$
Bayes Decision Rule: to minimize error
$$h(x) = \arg\max_{y \in \mathcal{Y}} [p(Y = y | X = x)]$$

**Discriminative models:** Use conditional distribution(?)
**MAP(Maximum A Posteriori Estimation)**
$$\hat{\theta} = \arg\max_{\theta} \log p(\theta|D) = \max_{\theta} \{\log p(D|\theta) + \log p(\theta)\}$$
**Generative Models:**
1. $p(X = x|Y = y)$: 给定类，输入的分布
2. $p(Y = y)$: 类自身的分布
Then $p(Y = y|X = x) \propto p(X = x|Y = y)p(Y = y)$
假设 1 是伯努利分布→Naïve Bayes
假设 1 是高斯分布→Gaussian Discriminant Analysis
均假设 2 是伯努利分布

---

### EM Algorithm:

$z$ : latent variables, $x$ :observed variables, $\theta$: parametric model, $p(x, z|\theta)$
**Key Idea:** What if we know $z$'s distribution $q(z)$?
**ELBO:** $L(q, \theta) = \sum_z q(z) \log\left(\frac{p(x, z|\theta)}{q(z)}\right)$
**KL-Div:** $KL[q(z)||p(z)] = \sum_z q(z) \log\left(\frac{q(z)}{p(z)}\right)$
**Properties:** $KL(p||q) \geq 0, KL(p||p) = 0$
$\log p(x|\theta) = L(q, \theta) + KL[q(z)||p(z|x, \theta)]$
**E-Step:** Maximizing $L(q, \theta)$ Over $q$ for Fixed $\theta = \theta^{old}$ best $q$ is $q^*(z) = p(z|x, \theta^{old}) \propto p(x|z, \theta^{old})q(z)$
**M-Step:** Maximizing $L(q, \theta)$ Over $\theta$ for Fixed $q$, equivalent to **maximize the expected complete data likelihood:** $\max_\theta \sum_z q(z) \log p(x, z|\theta)$ [Using SGD]
**MAP:** M-Step $\theta^{new} = \arg\max_\theta [L(q^*, \theta) + \log p(\theta)]$

**GMM:** $k$ clusters, 每个都是高斯分布,数据没有标签
**Dataset generation:** 1. random choose a cluster $z$, 2. random choose a point $x$ from cluster $z$
**EM approach:** E-step evaluate the responsibilities using current parameters for each point $\gamma_i^j$ M-step, use MLE to maximize total "Expectation" Match.

**Monte Carlo Estimator :** draw i.i.d. sample $\mathcal{D} = \{x_1, ..., x_n\}$ from distribution $p$, then use $\frac{1}{n}\sum_{i=1}^{n} f(x_i)$ to estimate $\mathbb{E}_p(f(x))$ [Unbiased, Variance $\frac{1}{n} Var_p(f(x))$]
**EM with Monte Carlo:**
M-step: $\theta^{new} = \cdots = \arg\max_\theta \mathbb{E}_{z \sim q} \log p(x, z|\theta) \approx \arg\max_\theta \frac{1}{T} \sum_{i=1}^{T} \log p(x, z_i|\theta)$ , $z$ is sampling from $q^* = p(z|x, \theta^{old})$
**Concrete sampling Method: 1.** Use the reverse of CDF **2.** Rejection Sampling: for $kq(x) > p(x)$, draw $x \sim q, u \sim U[0, kq(x)]$; reject $x$ if $u > p(x)$ [Problem: high rejection rate with high dimension]
**Markov Chain Monte Carlo:** Sampling from $p(x = j) = \pi_j \propto b_j$ ,where $b_j$ is computable, while $\sum b_j = B$ is unknown.
**Core idea:** Construct a Markov Chain with stationary distribution $\pi_j$.
**Method(Metropolis-Hastings):** Random transition matrix(proposal distribution) $Q: Q_{ij} = \mathcal{N}(x_j|x_i, \epsilon^2)$
1. start state $x_0$; 2. sample $x_*$ from $Q$; 3. $u \sim U[0,1]$, if $u > a_{ij} = \min\left(1, \frac{\pi_j Q_{ji}}{\pi_i Q_{ij}}\right)$, reject $x_*$; else, accept $x_*$.
**Detailed Balance cond.:** $\pi$ is stationary distribution $\Leftrightarrow \pi_i P_{ij} = \pi_j P_{ji}$.
· Monte Carlo with Markov Chain: $\frac{1}{T-T_0} \sum_{t=T_0}^{T-1} f(x_t) \to \mathbb{E}_{x \sim \pi} f(x)$
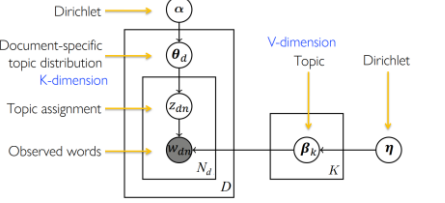
**Gibbs Sampling:** No rejection, with conditional distribution, a special case of Metro-Hast Algo with accept rate 1 (Proof TBD)
1. $x_1^{(t+1)} \sim p\left(x_1|x_2^{(t)}, x_3^{(t)}, ..., x_n^{(t)}\right)$
2. $x_2^{(t+1)} \sim p\left(x_2|x_1^{(t+1)}, x_3^{(t)}, ..., x_n^{(t)}\right)$
3. ...
4. $x_j^{(t+1)} \sim p\left(x_j|x_1^{(t+1)}, ..., x_{j-1}^{(t+1)}, x_{j+1}^{(t)}, ..., x_n^{(t)}\right)$
5. ...
6. $x_n^{(t+1)} \sim p\left(x_n|x_1^{(t+1)}, ..., x_{n-1}^{(t+1)}\right)$

---

Gaussian Discriminant Analysis的参数[数据有标签]
$$\phi = \frac{\sum_{i=1}^{n} 1\{y_i = +1\}}{n}, \mu_+ = \frac{\sum_{i=1}^{n} 1\{y_i = +1\}x_i}{\sum_{i=1}^{n} 1\{y_i = +1\}}, \mu_- = \frac{\sum_{i=1}^{n} 1\{y_i = -1\}x_i}{\sum_{i=1}^{n} 1\{y_i = -1\}}$$
$$\Sigma = \frac{1}{n} \sum_{i=1}^{n} (x_i - \mu_{y_i})(x_i - \mu_{y_i})^T$$

**Dirichlet Distribution:** $p(\theta|\alpha) = \frac{1}{B(\alpha)} \prod_{k=1}^{K} \theta_k^{\alpha_k - 1}$,
$\theta$ must satisfy $\sum_{i=1}^{K} \theta_i = 1$
Explanation: $\alpha_k$ is the importance of class $k$
**Bayesian Estimation:** The expectation of parameter using posterior as distribution
$$\hat{\theta} = \mathbb{E}_{\theta \sim p(\theta|D)}[\theta] = \int \theta \cdot p(\theta|D)d\theta$$

**Dirichlet-Multinomial Conjugate**
$\beta \sim Dir(\beta|\eta), w \sim Mult(1, \beta), n$ is count vector
The posterior Distribution is also Dirichlet! ➔ Simplify Below!
$$p(\beta|\eta, \mathcal{W}) = Dir(\beta|\eta + n)$$
**MAP:** $\hat{\beta}_{MAP} = \max_\beta Dir(\beta, \eta + n), (\hat{\beta}_{MAP})_i = \frac{\eta_i + n_i - 1}{\sum_{j=1}^{V} \eta_j + N - V}$
**Bayes Est.:** $\hat{\beta}_{Bayes} = \int \beta Dir(\beta|\eta + n)d\beta, (\hat{\beta}_{Bayes})_i = \frac{\eta_i + n_i}{\sum_{j=1}^{V} \eta_j + N}$

### Hidden Markov Model

Hidden Variable $z_t$: Markov Chain, Transition $P(z_t|z_{t-1})$
Observation $x_t$: only depend on $z_t$, Emission $P(x_t|z_t)$
**Inference for HMM:** know $x$, guess $z$
**Sequence error:** $L(Z, \hat{Z}) = 1(\hat{Z} \neq Z), \hat{Z} = \arg\max_Z p(Z|X) \propto \arg\max_Z p(Z, X)$ [use joint distribution]
**Viterbi Algorithm:** negative log-probability, change to a shortest path model (TBD)
**State error:** $L(Z, \hat{Z}) = \sum_t 1(\hat{z}_t \neq z_t), \hat{z}_t = \max_{z_t} p(z_t|X)$ [use marginal distribution]
**Forward-Backward Algorithm,** $p(z_t|X) \propto P(z_t|x_{1:t})P(x_{t+1:T}|z_t) = \alpha_t(z_t)\beta_t(z_t)$
**Forward part:** $\alpha_t(z_t) = P(z_t|x_{1:t})$
1. Push by transition $P(z_t|x_{1:t-1}) = \sum_{z_{t-1}} P(z_t|z_{t-1})\alpha_{t-1}(z_{t-1})$
2. Reweighted by emission $\alpha_t(z_t) = P(z_t|x_{1:t}) \propto P(x_t, z_t|x_{1:t-1}) = P(x_t|z_t)P(z_t|x_{1:t-1})$
**Backward part:** $\beta_t(z_t) = P(x_{t+1:T}|z_t)$
$$\beta_t(z_t) = \sum_{z_{t+1}} \beta_{t+1}(z_{t+1})P(x_{t+1}|z_{t+1})P(z_{t+1}|z_t)$$