
No Coding Farmer

Tao Hu

Department of Computer Science
Peking University
No.5 Yiheyuan Road Haidian District, Beijing, P.R.China
taohu@pku.edu.cn

Abstract

Some Miscellaneous Summary.

Contents

1	Expectation Maximization Introduction	3
1.1	EM Induction	3
1.2	EM convergence proof	3
1.3	Different Writing Style of EM Algorithm	3
2	EM applications	4
2.1	Gaussian Mix Model	4
2.2	Hidden Markov Model	4
2.3	Naive Bayesian	5
2.4	other papers	5
3	VAE	5
4	ADMM	5
5	Key steps you must know when building a DL Framework	8
5.1	Convolution	8
6	R-PCA	9
6.1	Solve RPCA by ADMM	9
6.2	Adaptive Penalty for ADMM	9
7	SFM	10
8	Reinforcement Learning	10
8.1	Model-based Method	12
8.1.1	Policy Iteration	12
8.1.2	Value Iteration	14
8.2	Model-Free Method	14
8.2.1	on-policy TD	14
8.2.2	on-policy MCMC	15
8.2.3	off-policy method	15
8.2.4	Comparisons: DP, MC, TD	15

1 Expectation Maximization Introduction

1.1 EM Induction

$$L(\theta) = \sum_{i=1}^M \log p(X; \theta) = \sum_{i=1}^M \log \sum_z p(X, Z; \theta)$$

let θ_i be some distribution over z 's ($\sum_z \theta_i(z) = 1, \theta_i(z) \geq 0$)

$$\begin{aligned} & \sum_i \log p(X^{(i)}; \theta) \\ &= \sum_i \log \sum_{Z^{(i)}} \theta_i(Z^{(i)}) \frac{p(X^{(i)}, Z^{(i)}; \theta)}{\theta_i(Z^{(i)})} \\ &\geq \sum_i \sum_{Z^{(i)}} \theta_i(Z^{(i)}) \log \frac{p(X^{(i)}, Z^{(i)}; \theta)}{\theta_i(Z^{(i)})} \quad (f(x) = \log x \text{ is concave.}) \end{aligned}$$

$$\text{let } \frac{p(X^{(i)}, Z^{(i)}; \theta)}{\theta_i(Z^{(i)})} = C$$

the equality can be only reached when $\frac{p(X^{(i)}, Z^{(i)}; \theta)}{\theta_i(Z^{(i)})}$ is a constant.

we can get: $\sum_i \frac{p(X^{(i)}, Z^{(i)}; \theta)}{C} = 1$ namely: $\sum_i p(X^{(i)}, Z^{(i)}; \theta) = C$

further induction: $\theta_i(Z^{(i)}) = \frac{p(X^{(i)}, Z^{(i)}; \theta)}{\sum_i p(X^{(i)}, Z^{(i)}; \theta)} = p(Z^{(i)} | X^{(i)}; \theta)$

so the procedure of EM algorithm is:

Repeat Until Convergence:

- E-step: for each i , get $Q_i(Z^{(i)}) = p(Z^{(i)} | X^{(i)}; \theta)$
- M-step: $\theta := \argmax_{\theta} \sum_i \sum_{Z^{(i)}} Q_i(Z^{(i)}) \log \frac{p(X^{(i)}, Z^{(i)}; \theta)}{Q_i(Z^{(i)})}$

1.2 EM convergence proof

$$\text{let } l(\theta^{(t)}) = \sum_i \sum_{Z^{(i)}} Q_i^{(t)}(Z^{(i)}) \log \frac{p(X^{(i)}, Z^{(i)}; \theta)}{Q_i^{(t)}(Z^{(i)})}$$

then, we have the following inequality:

$$\begin{aligned} & l(\theta^{(t+1)}) \\ &\geq \sum_i \sum_{Z^{(i)}} Q_i^{(t)}(Z^{(i)}) \log \frac{p(X^{(i)}, Z^{(i)}; \theta^{(t+1)})}{Q_i^{(t)}(Z^{(i)})} \\ &\geq \sum_i \sum_{Z^{(i)}} Q_i^{(t)}(Z^{(i)}) \log \frac{p(X^{(i)}, Z^{(i)}; \theta^{(t)})}{Q_i^{(t)}(Z^{(i)})} \\ &\geq l(\theta^{(t)}) \end{aligned}$$

the first inequality is because: $l(\theta) \geq \sum_i \sum_{Z^{(i)}} Q_i^{(t)}(Z^{(i)}) \log \frac{p(X^{(i)}, Z^{(i)}; \theta)}{Q_i^{(t)}(Z^{(i)})} \forall \theta, Q_i$

the second inequality is because of the maximum of the M-step.

Hence, EM causes the likelihood to converge monotonically.

1.3 Different Writing Style of EM Algorithm

There are many writing style of EM algorithm. here I just mention the book <Statistics Learning Method> by LiHang who is very famous in China.

EM algorithm from LiHang(Li-version):

Algorithm 1 EM from LIHang

Require: observation X , hidden variable Z , joint distribution $P(X, Z|\theta)$, conditional distribution $P(Z|Y, \theta)$
while Not convergence **do**
 E-Step: let $\theta^{(i)}$ is the i -th estimate of θ ,
 $Q(\theta, \theta^{(i)}) = E_z[\log P(X, Z|\theta)|X, \theta^{(i)}] = \sum_Z \log P(X, Z|\theta)P(Z|X, \theta^{(i)})$
 M-step: $\theta^{(i+1)} = \arg\max_{\theta} Q(\theta, \theta^{(i)})$
end while
output model parameter θ

it seems that Li-version is different from the above version. however, they are the same. because:

- the above version just consider every data, so that it include subscript i . however Li-version only consider one data.
- the above version can be transformed to Li-version.

$$\begin{aligned} & \sum_Z Q(Z) \log \frac{P(X, Z; \theta)}{Q(Z)} \\ &= \sum_Z P(Z|X; \theta^{(t)}) \log \frac{P(X, Z; \theta)}{P(Z|X; \theta^{(t)})} \\ &= \sum_Z P(Z|X; \theta^{(t)}) \log P(X, Z; \theta) - \sum_Z P(Z|X; \theta^{(t)}) \log P(Z|X; \theta^{(t)}) \end{aligned}$$

as the variable is θ , so $\sum_Z P(Z|X; \theta^{(t)}) \log P(Z|X; \theta^{(t)})$ can be removed.

- $Q(\theta, \theta^{(i)}) = \sum_Z \log P(X, Z|\theta)P(Z|X, \theta^{(i)})$ can be also written as $Q(\theta, \theta^{(i)}) = \sum_Z \log P(X, Z|\theta)P(Z, X, \theta^{(i)})$, because X is a observation.

2 EM applications

2.1 Gaussian Mix Model

GMM can be solved by EM. notice here we use the expectation of EM:

$$\begin{aligned} & Q(\theta, \theta^{(i)}) \\ &= E_{\gamma}[\log P(y, \gamma|\theta)|y, \theta^{(i)}] \\ &= E[\sum_{k=1}^K [n_k \log \alpha_k + \sum_{j=1}^N \gamma_{jk} [\log \frac{1}{\sqrt{2\pi}} - \log \sigma_k - \frac{1}{2\sigma_k^2} (y_j - \mu_k)^2]]] \\ &= \sum_{k=1}^K [(E\gamma_{jk}) \log \alpha_k + \sum_{j=1}^N (E\gamma_{jk}) [\log \frac{1}{\sqrt{2\pi}} - \log \sigma_k - \frac{1}{2\sigma_k^2} (y_j - \mu_k)^2]] \end{aligned}$$

here $(E\gamma_{jk})$ can be easily calculated.

$\hat{\mu}_k, \hat{\sigma}_k^2$ can be acquired by derivation.

$\hat{\alpha}_k$ can be acquired by the derivation on the Lagrangian ($\sum_i^K \alpha_k = 1$).

2.2 Hidden Markov Model

HMM Learning Method is also called Baum-Welch algorithm. the target is learning $\lambda = (A, B, \pi)$.

Q function is:

$$\begin{aligned} Q(\lambda, \bar{\lambda}) &= \sum_I \log P(O, I|\lambda) P(O, I|\bar{\lambda}) \\ P(O, I, \lambda) &= \pi_{i_1} b_{i_1}(o_1) a_{i_1 i_2} b_{i_2}(o_2) \dots a_{i_{T-1} i_T} b_{i_T}(o_T) \end{aligned}$$

so the Q function can also be written as:

$$Q(\lambda, \bar{\lambda}) = \sum_I \log \pi_{i1} P(O, I | \bar{\lambda}) + \sum_I (\sum_{t=1}^{T-1} \log a_{i,i+1}) P(O, I | \bar{\lambda}) + \sum_I (\sum_{t=1}^T \log b_{it}(o_t)) P(O, I | \bar{\lambda})$$

note here: I is not only one state. it includes state length from 1 to T, which all start from i_1

so we can solve the maximum of Q function by derivation on the Lagrangian polynomial (because exists these limitations: $\sum_{i=1}^N \pi_i = 1, \sum_{j=1}^N a_{ij} = 1, \sum_{i=1}^M b_i = 1$)

2.3 Naive Bayesian

2.4 other papers

We can use softmax to model transition probability, normal distribution to model emission probability.

it's a good example in Car that Knows Before You Do: Anticipating Maneuvers via Learning Temporal Driving Models, the AIO-HMM can be more complicated, which can be enriched by the graphic model by M.I Jordon.

3 VAE

here is a complete VAE tutorial [2]

$$\begin{aligned} \max \quad & \log P(x) \\ \text{lhs} = & \log \int P(x, z) dz \\ = & \log \int P(x/z) p(z) dz \\ = & \log \int \frac{P(x/z)}{q(z/x)} q(z/x) p(z) dz \\ = & \log E_{q(z/x)} \left[\frac{P(x/z)}{q(z/x)} p(z) \right] \\ \text{jensen's inequality, we can know: } & \geq E_{q(z/x)} [\log \frac{P(x/z)}{q(z/x)} p(z)] \\ = & E_{q(z/x)} [\log p(x/z)] + E_{q(z/x)} [\log \frac{p(z)}{q(z/x)}] \\ = & E_{q(z/x)} [\log p(x/z)] - E_{q(z/x)} [\log \frac{q(z/x)}{p(z)}] \\ = & E_{q(z/x)} [\log p(x/z)] - KL(q(z/x) || p(z)) \end{aligned}$$

4 ADMM

minimize $H(u) + G(v)$
subject to $Au + Bv = b$

$$\max_{\lambda} \min_{u,v} H(u) + G(v) + \langle \lambda, b - Au - Bv \rangle + \frac{\tau}{2} \|b - Au - Bv\|^2$$

Alternating Direction Method of Multipliers

$$\begin{aligned} u_{k+1} &= \arg \min_u H(u) + \langle \lambda_k, -Au \rangle + \frac{\tau}{2} \|b - Au - Bv_k\|^2 \\ v_{k+1} &= \arg \min_v G(v) + \langle \lambda_k, -Bv \rangle + \frac{\tau}{2} \|b - Au_{k+1} - Bv\|^2 \\ \lambda_{k+1} &= \lambda_k + \tau(b - Au_{k+1} - Bv_{k+1}) \end{aligned}$$

Distributed Problems

minimize $g(x) + \sum_i f_i(x)$
example: sparse least squares:
minimize $\mu \|x\|_1 + \frac{1}{2} \|Ax - b\|^2$

$$\text{minimize } \mu \|x\|_1 + \sum_i \frac{1}{2} \|A_i x - b_i\|^2$$

data stored on different servers

Transpose Reduction

minimize $\frac{1}{2} \|Ax - b\|^2$

$$x^* = (A^T A)^{-1} A^T b$$

distributed computation:

$$A^T b = \sum A_i^T b_i$$

$$A^T A = \sum A_i^T A_i$$

Unwrapped ADMM

minimize $g(x) + f(Ax) = g(x) + \sum_i f_i(A_i x)$

Example: SVM

minimize $\frac{1}{2} \|x\|^2 + h(Ax)$

A = data, h = hinge loss

Unwrapped form

minimize $\frac{1}{2} \|x\|^2 + h(z)$

subject to $z = Ax$

Transpose Reduction ADMM

scaled augmented Lagrangian:

minimize $\frac{1}{2} \|x\|^2 + h(z) + \frac{\tau}{2} \|z - Ax - \lambda\|^2$

ADMM:

$$x^{k+1} = \min_x \frac{1}{2} \|x\|^2 + \frac{\tau}{2} \|z^k - Ax + \lambda^k\|^2$$

$$z^{k+1} = \min_z h(z) + \frac{\tau}{2} \|z - Ax^{k+1} + \lambda^k\|^2$$

$$\lambda^{k+1} = \lambda^k + z^{k+1} - Ax^{k+1}$$

Minimization Steps

minimize $l(a_3) + \frac{1}{2} \|z_2 - W_1 a_1\|^2 + \frac{1}{2} \|a_2 - \sigma(z_2)\|^2 + \frac{1}{2} \|z_3 - W_2 a_2\|^2 + \frac{1}{2} \|a_3 - \sigma(z_3)\|^2$

Solve for weight: least squares(convex)

Solve for activations: least squares + ridge penalty(convex)

Solve for inputs: coordinate-minimization (non-convex but global)

Lagrange Multipliers

minimize $l(a_3) + \frac{1}{2} \|z_2 - W_1 a_1\|^2 + \frac{1}{2} \|a_2 - \sigma(z_2)\|^2 + \langle \lambda_1, z_2 - W_1 a_1 \rangle + \langle \lambda_2, a_2 - \sigma(z_2) \rangle + \frac{1}{2} \|z_3 - W_2 a_2\|^2 + \frac{1}{2} \|a_3 - \sigma(z_3)\|^2 + \langle \lambda_3, z_3 - W_2 a_2 \rangle + \langle \lambda_4, a_3 - \sigma(z_3) \rangle$

unstable because of non-linear constraints

Bregman Iteration

minimize $l(a_3) + \langle \lambda, a_3 \rangle + \frac{1}{2} \|z_2 - W_1 a_1\|^2 + \frac{1}{2} \|a_2 - \sigma(z_2)\|^2 + \frac{1}{2} \|z_3 - W_2 a_2\|^2 + \frac{1}{2} \|a_3 - \sigma(z_3)\|^2$

HOG feature dimension: 648

mid layer 1 num: 100

mid layer 2 num: 50

output layer: 1

Algorithm 2 ADMM_NN

Inputs:

data number: $n=10000$,
data dimension: $m=648$,
hidden layer 1 unit number: $a=100$
hidden layer 2 unit number: $b=50$
output layer unit number: 1
 a_0 m-n dimension,
 W_1 : a-m dimension
 z_1 : a-n dimension
 a_1 : a-n dimension
 W_2 : b-a dimension
 z_2 : b-n dimension
 a_2 : b-n dimension
 W_3 : 1-b dimension
 z_3 : 1-n dimension
labels: y 1-n dimension
 λ : 1-n dimension
activation function h is ReLu.

Initialize:

allocate $\{a_l\}_{l=1}^L, \{z_l\}_{l=1}^L$ with i.i.d Gaussian Distribution, and λ

Cache: a_0^\dagger

Warm Start:

for $i=1, \dots, 100$ **do**

for $l=1, 2, \dots, L-1$ **do**

$W_l \leftarrow z_l a_{l-1}^\dagger$

$a_l \leftarrow (\beta_{l+1} W_{l+1}^T W_{l+1} + \gamma_l I)^{-1} (\beta_{l+1} W_{l+1}^T z_{l+1} + \gamma_l h_l(z_l))$

$z_l \leftarrow \operatorname{argmin}_z \gamma_l \|a_l - h_l(z)\|^2 + \beta_l \|z - W_l a_{l-1}\|^2$

end for

$W_L \leftarrow z_L a_{L-1}^\dagger$

$z_L \leftarrow \operatorname{argmin}_z l(z, y) + \langle z, \lambda \rangle + \beta_L \|z - W_L a_{L-1}\|^2$

end for

Start ADMM:

while not converge **do**

for $l=1, 2, \dots, L-1$

do $W_l \leftarrow z_l a_{l-1}^\dagger$

$a_l \leftarrow (\beta_{l+1} W_{l+1}^T W_{l+1} + \gamma_l I)^{-1} (\beta_{l+1} W_{l+1}^T z_{l+1} + \gamma_l h_l(z_l))$

$z_l \leftarrow \operatorname{argmin}_z \gamma_l \|a_l - h_l(z)\|^2 + \beta_l \|z - W_l a_{l-1}\|^2$

end for

$W_L \leftarrow z_L a_{L-1}^\dagger$

$z_L \leftarrow \operatorname{argmin}_z l(z, y) + \langle z, \lambda \rangle + \beta_L \|z - W_L a_{L-1}\|^2$

$\lambda \leftarrow \lambda + \beta_L (z_L - W_L a_{L-1})$

end while

z_l argmin procedure:

$$z_l = \begin{cases} \max(\frac{a_l \gamma_l + W_l a_{l-1} \beta_l}{\gamma_l + \beta_l}, 0) & z \geq 0 \\ \min(W_l a_{l-1}, 0) & z \leq 0 \end{cases}$$

choose one minimizer z from two choices.

z_L argmin procedure:

when $y_i = 0$:

$$f(z) = \beta z^2 - (2\beta w_a - \lambda)z + \max(z, 0)$$

$$z^* = \max(\frac{2\beta w_a - \lambda - 1}{2\beta}, 0) \text{ or}$$

$$z^* = \min(\frac{2\beta w_a - \lambda}{2\beta}, 0)$$

choose one which make $f(z)$ smaller.

when $y_i = 1$:

$$f(z) = \beta z^2 - (2\beta w_a - \lambda)z + \max(1 - z, 0)$$

$$z^* = \max(\frac{2\beta w_a - \lambda}{2\beta}, 1) \text{ or}$$

$$z^* = \min(\frac{2\beta w_a - \lambda + 1}{2\beta}, 1)$$

choose one which make $f(z)$ smaller.

z_L argmin procedure(when l is a standard hinge loss):

$$\text{when } y_i = -1: f(z) = \max(1 + z, 0) + \lambda z + \beta(z^2 - 2w_a z)$$

$$z^* = \min(\frac{2\beta w_a - \lambda}{2\beta}, -1) \text{ or}$$

$$z^* = \max(\frac{2\beta w_a - \lambda - 1}{2\beta}, -1) \text{ choose one which make } f(z) \text{ smaller.}$$

when $y_i = 1$:

$$f(z) = \max(1 - z, 0) + \lambda z + \beta(z^2 - 2w_a z)$$

$$z^* = \min(\frac{2\beta w_a - \lambda + 1}{2\beta}, 1) \text{ or}$$

$$z^* = \max(\frac{2\beta w_a - \lambda - 1}{2\beta}, 1)$$

choose one which make $f(z)$ smaller.

5 Key steps you must know when building a DL Framework

5.1 Convolution

convert the convolution to matrix multiplication like the fully-connected network.

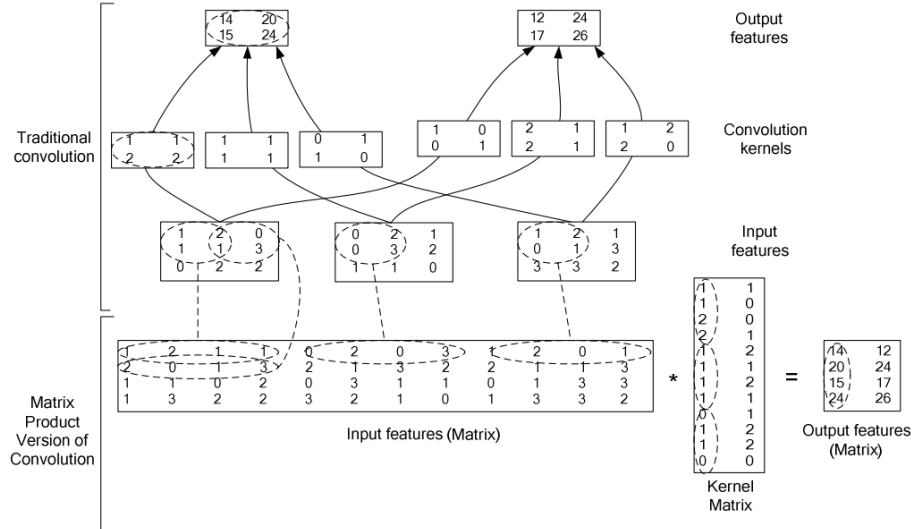


Figure 2. Example convolution operations in a convolutional layer (biases, sub-sampling, and non-linearity omitted). The top figure presents the traditional convolution operations, while the bottom figure presents the matrix version.

figure is from [1].

Let's notate:

H:image height

W:image width

in: input image number

out: output image number

K: convolution kernel size

the matrix product version of convolution, the dimension of two matrix is:

$$(H*W) * (in*K*K)$$

$$(in*K*K) * out$$

the operation above is like matlab function: `img2col`
`im2col(A,[m n],block_type), where block_type="sliding".`

GEMM(GEneral Matrix to Matrix Multiplication) is at the heart of deep learning. <https://petewarden.com/2015/04/20/why-gemm-is-at-the-heart-of-deep-learning/>

6 R-PCA

RPCA problem:

$$\min_{A,E} \|A\|_* + \lambda \|E\|_1$$

$$S.t \ D=A+E$$

RPCA dual problem:

Augmented Lagrangian is :

$$\begin{aligned} A_t(A, E; \Lambda) &= \min_{A,E} EL(A, E; \Lambda) \\ &= \min_{A,E} \|A\|_* + \Lambda \|E\|_1 + \langle \Lambda, D - A - E \rangle \\ &= \min_A \|A\|_* - \langle \Lambda, A \rangle + \min_E \lambda \|E\|_1 - \langle \Lambda, E \rangle + \langle \Lambda, D \rangle \end{aligned}$$

both of the sub-problem is conjugate function, according to the property of conjugate function :

$$\begin{aligned} A_t(A, E; \Lambda) &= \langle \Lambda, D \rangle \\ S.t \quad &\|\Lambda\|_2 \leq 1, \|\Lambda\|_\infty \leq \lambda \end{aligned}$$

so the dual problem is:

$$\begin{aligned} \max_{\Lambda} \quad &\langle \Lambda, D \rangle \\ S.t \quad &\|\Lambda\|_2 \leq 1, \|\Lambda\|_\infty \leq \lambda \end{aligned}$$

6.1 Solve RPCA by ADMM

the ADMM sub-problem is:

A-sub-problem:

$$A_{k+1} = \operatorname{argmin}_A \|A\|_* + \frac{\beta}{2} \|D - A - E_k + \Lambda_k / \beta\|_F^2$$

E-sub-problem:

$$E_{k+1} = \operatorname{argmin}_E \lambda \|E\|_1 + \frac{\beta}{2} \|D - A_{k+1} - E + \Lambda_k / \beta\|_F^2$$

E-sub-problem has closed-form solution as follows:

$$E_{k+1} = S_{\lambda/\beta}(D - A_{k+1} + \Lambda_k / \beta).$$

$S_\epsilon = \operatorname{sgn}(x) \max(|x| - \epsilon, 0)$, which is the same form as shrinkage.

A-sub-problem has a closed-form solution offered by Singular Value Thresholding(SVT): suppose that the SVD of $W = D - E_k + \Lambda_k / \beta$ is $W = U \Sigma V^T$, then the optimal solution is $A = U S_{\beta^{-1}}(\Sigma) V^T$.

6.2 Adaptive Penalty for ADMM

Lin et al.[3] suggest updating the penalty parameter β as follows:

$$\beta_{k+1} = \min(\beta_{max}, \rho \beta_k)$$

where ρ_{max} is an upper bound of $\{\beta_k\}$. the value of ρ is defined as:

$$\rho = \begin{cases} \rho_0 & \text{if } \frac{\beta_k \max(\sqrt{\eta_A} \|x_{k+1} - x_k\|_2, \sqrt{\eta_B} \|y_{k+1} - y_k\|_2)}{\|c\|_2} < \epsilon_2 \\ 1 & \text{otherwise} \end{cases}$$

where η_A, η_B is linearized Taylor second-order factor.

7 SFM

<https://www.robots.ox.ac.uk/~vgg/hzbook/hzbook2/HZepipolar.pdf>

- F is a rank 2 homogeneous matrix with 7 degrees of freedom.
- **Point correspondence:** If \mathbf{x} and \mathbf{x}' are corresponding image points, then $\mathbf{x}'^T F \mathbf{x} = 0$.
- **Epipolar lines:**
 - ◊ $\mathbf{l}' = F \mathbf{x}$ is the epipolar line corresponding to \mathbf{x} .
 - ◊ $\mathbf{l} = F^T \mathbf{x}'$ is the epipolar line corresponding to \mathbf{x}' .
- **Epipoles:**
 - ◊ $F \mathbf{e} = \mathbf{0}$.
 - ◊ $F^T \mathbf{e}' = \mathbf{0}$.
- **Computation from camera matrices P, P' :**
 - ◊ General cameras, $F = [\mathbf{e}']_{\times} P' P^+$, where P^+ is the pseudo-inverse of P , and $\mathbf{e}' = P' \mathbf{C}$, with $P \mathbf{C} = \mathbf{0}$.
 - ◊ Canonical cameras, $P = [\mathbf{I} \mid \mathbf{0}]$, $P' = [\mathbf{M} \mid \mathbf{m}]$, $F = [\mathbf{e}']_{\times} \mathbf{M} = \mathbf{M}^{-T} [\mathbf{e}]_{\times}$, where $\mathbf{e}' = \mathbf{m}$ and $\mathbf{e} = \mathbf{M}^{-1} \mathbf{m}$.
 - ◊ Cameras not at infinity $P = K[\mathbf{I} \mid \mathbf{0}]$, $P' = K'[\mathbf{R} \mid \mathbf{t}]$, $F = K'^{-T} [\mathbf{t}]_{\times} \mathbf{R} K^{-1} = [\mathbf{K}' \mathbf{t}]_{\times} K' \mathbf{R} K^{-1} = K'^{-T} \mathbf{R} K^T [\mathbf{K} \mathbf{R}^T \mathbf{t}]_{\times}$.

Figure 1: Summary of Fundamental matrix properties

8 Reinforcement Learning

Markov Decision Process (**MDP**) is tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma)$:

- r is a reward function, $r(s, a, s')$
- γ is a discount factor
- \mathcal{P} is the transition probability distribution:
probability from state s with action a to state $s' : P(s' | s, a)$
- \mathcal{S} is a finite set of states.
- \mathcal{A} is a finite set of actions.

Markov Property:

$$P(s_{t+1} | s_t) = P(s_{t+1} | s_1, \dots, s_t)$$

Stochastic Policy:

$\pi(a|s) = P(a_s = a|s_t = s)$, $\pi(a|s)$ here is a probability, we can often see $\pi(s)$, which returns a, means under policy π and state s , you should better take action a .

Value function:

$$G_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

State-value function:

$$V_{\pi}(s) = E_{\pi}(G_t|s_t = s)$$

Action-value function:

$$Q_{\pi}(s, a) = E_{\pi}(G_t|s_t = s, a_t = a)$$

Bellman Equation:

$$\begin{aligned} V_{\pi}(s) &= E_{\pi}(r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots | s_t = s) \\ &= E_{\pi}(r_{t+1} + \gamma(r_{t+2} + \gamma r_{t+3} + \dots) | s_t = s) \\ &= E_{\pi}(r_{t+1} + \gamma G_{t+1} | s_t = s) \\ &= E_{\pi}(r_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s) \end{aligned}$$

For state-value function, Bellman Equation can be written as:

$$\begin{aligned} V_{\pi}(s) &= E_{\pi}(r_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s) \\ &= \sum_{a \in A} \pi(a|s) \sum_{s' \in S} P(s'|s, a) [r(s, a, s') + \gamma V_{\pi}(s')] \\ r(s, a, s') &\text{ is same with } r_{t+1} \text{ to some extent.} \end{aligned}$$

For action-value function, Bellman Equation can be written as:

$$\begin{aligned} Q_{\pi}(s, a) &= E_{\pi}(r(s, a, s') + \gamma Q_{\pi}(s', a') | s, a) \\ &= \sum_{s' \in S} P(s'|s, a) [r(s, a, s') + \gamma \sum_{a' \in A} \pi(a', s') Q_{\pi}(s', a')] \end{aligned}$$

Normally, we just assume the $\pi(a|s)$, $p(s'|s, a)$, $r(s, a, s')$ are known (namely the MDP is known). so we can solve the linear equation above. however, when data become huge, it is not feasible to solve the Bellman Equation directly.

optimal value function:

the optimal state-value function $V_*(s)$ is :

$$V_*(s) = \max_{\pi} V_{\pi}(s)$$

the optimal action-value function $Q_*(s, a)$ is:

$$Q_*(s, a) = \max_{\pi} Q_{\pi}(s, a)$$

$$Q^*(s, a) = \max_{\pi} Q_{\pi}(s, a)$$

following important properties:

$$\begin{aligned} Q^*(s, a) &= E[r_{t+1} + \gamma V^*(s_{t+1}) | s_t = s, a_t = a] \\ V^*(s) &= \max_{a \in A} Q^*(s, a) \\ Q^*(s, a) &= \sum_{s' \in S} p(s'|s, a) [r(s, a, s') + \gamma V^*(s')] \end{aligned}$$

The goal for any MDP is finding the optimal value function.

Or equivalent an optimal policy π^* for any policy π , $V_{\pi^*}(s) \geq V_{\pi}(s)$, $\forall s \in S$

how do we take action after we obtain the optimal $V^*(s)$? if we known the $p(s'|s, a)$, $r(s, a, s')$, then we can get: $Q^*(s, a) = \sum_{s' \in S} p(s'|s, a) [r(s, a, s') + \gamma V^*(s')]$. after $Q^*(s, a)$ is acquired, it's more trivial to take action a when it's at state s .

Relation between Q and V Functions

reference: <http://www.cs.cmu.edu/~mgormley/courses/10601-s17/slides/lecture26-ri.pdf>

Q from V:

$$Q^\pi(s, a) = \sum_{s' \in \mathcal{S}} P(s' | s, a) [R(s, a, s') + \gamma V^\pi(s')]$$

V from Q:

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a | s) Q^\pi(s, a)$$

V and Q

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} p(s' | s, a) [r(s, a, s') + \gamma V^\pi(s')]$$

Q and V

$$Q^\pi(s, a) = \sum_{s' \in \mathcal{S}} p(s' | s, a) [r(s, a, s') + \gamma \sum_{a' \in \mathcal{A}} \pi(a' | s') Q^\pi(s', a')]$$

more complicated..

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} p(s' | s, a) [r(s, a, s') + \gamma \sum_{a' \in \mathcal{A}} \pi(a' | s') Q^\pi(s', a')]$$

mainly can be classified into two main approaches:

Model based approaches:

First we will discuss methods that need to know the model:

$$P(s' | s, a) \text{ and } R(s, a, s').$$

☐ **Policy Iteration**

☐ **Value Iteration**

Model-free approaches:

Then we will discuss “model-free” methods that do NOT need to know the model: $P(s' | s, a)$ and $R(s, a, s')$.

☐ **Monte Carlo Method**

☐ **TD Learning**

34

Figure 2: Structure of RL

8.1 Model-based Method

8.1.1 Policy Iteration

One drawback of policy iteration is that each iteration involves policy evaluation.

1. Initialization

$V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ for all $s \in \mathcal{S}$.
 $\pi(s)$ is a deterministic policy.
 $\delta > 0$ is a small threshold parameter.

2. Policy Evaluation

```
repeat
   $\Delta \leftarrow 0$ 
  for all  $s \in \mathcal{S}$  do:
     $v \leftarrow V(s)$ 
     $a \leftarrow \pi(s)$ 
     $V(s) \leftarrow \sum_{s' \in \mathcal{S}} P(s'|s, a)[R(s, a, s') + \gamma V(s')]$ 
     $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
  end for
until  $\Delta < \delta$ 
```

3. Policy Improvement

```
policyStable  $\leftarrow$  true
for all  $s \in \mathcal{S}$  do:
   $b \leftarrow \pi(s)$ 
   $\pi(s) \leftarrow \arg \max_{a \in \mathcal{A}(s)} \sum_{s' \in \mathcal{S}} P(s'|s, a)[R(s, a, s') + \gamma V(s')]$ 
  if  $b \neq \pi(s)$  then
    policyStable  $\leftarrow$  false
  end if
end for
if policyStable then
  STOP
else
  Go to 2 (Policy Evaluation)
end if
```

Policy Improvement just improve the policy $\pi(s)$, then when we back to policy evaluation, the next action a is determined by the policy $\pi(s)$ which was updated in policy improvement. thus the relationship of policy improvement and policy evaluation is founded.

8.1.2 Value Iteration

Value Iteration

Main idea:

Use the Bellman equation of V^* instead of V^π

The greedy operator:

$$[T^*V](s) := \max_{a \in \mathcal{A}(s)} \sum_{s' \in \mathcal{S}} P(s'|s, a) [R(s, a, s') + \gamma V(s')]$$

V^* is the solution of $V = T^*(V)$ fixpoint iteration.

The value iteration update:

$k = 0$ and $V_0(s) \in \mathbb{R}$ for all $s \in \mathcal{S}$

repeat

for all $s \in \mathcal{S}$ **do:**

$$V_{k+1}(s) \leftarrow \max_{a \in \mathcal{A}(s)} \sum_{s' \in \mathcal{S}} P(s'|s, a) [R(s, a, s') + \gamma V_k(s')]$$

end for

$k \leftarrow k + 1$

until $V_k(\cdot)$ converged

37

- 1, value iteration converges to the true solution of Bellman optimal equations
- 2, Learn optimal value function directly, unlike policy iteration, there is no explicit policy.

8.2 Model-Free Method

8.2.1 on-policy TD

TD0

TD(n)

TD λ

8.2.2 on-policy MCMC

8.2.3 off-policy method

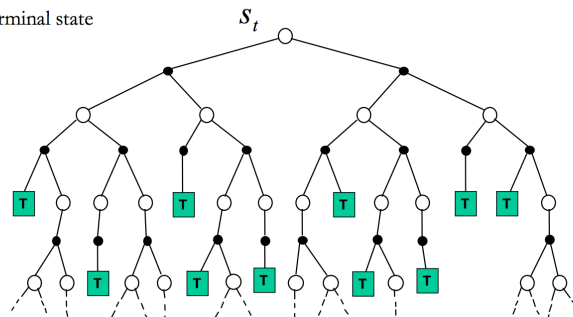
8.2.4 Comparisons: DP, MC, TD

Comparisons: DP, MC, TD

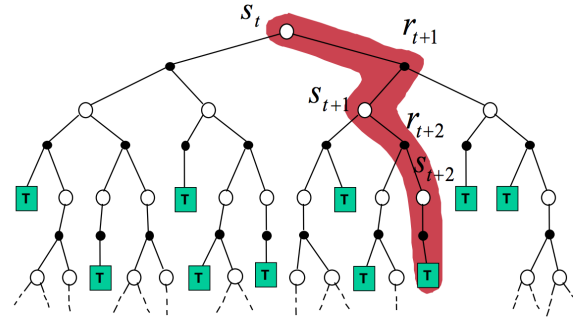
- They all estimate V^π
- DP: $V_k(s_t) \approx E_\pi(r_{t+1} + \gamma V_{k-1}(s_{t+1}) | s_t)$
 - Estimate comes from the Bellman equation
 - It needs to know the model
- TD: $V_k(s_t) \approx (r_{t+1} + \gamma V_{k-1}(s_{t+1}))$
 - Expectation is approximated with random samples
 - Doesn't need to wait for the end of the episodes.
- MC: $V_k(s_t) \approx R_t(s_t)$
 - Expectation is approximated with random samples
 - It needs to wait for the end of the episodes

MDP Backup Diagrams

- White circle: state
- Black circle: action
- T: terminal state

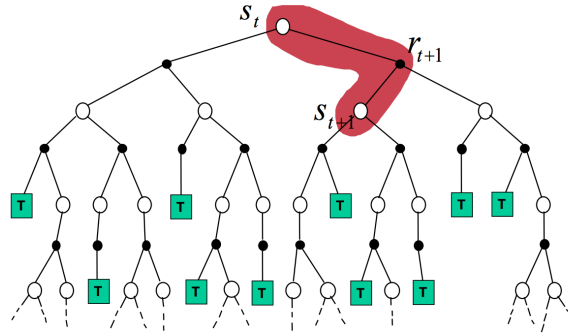


Monte Carlo Backup Diagram



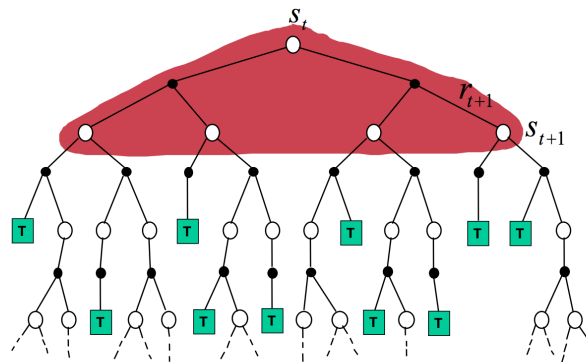
MC estimate: $V_k(s_t) := V_{k-1}(s_t) + \alpha_k \cdot (R_k(s_t) - V_{k-1}(s_t))$

Temporal Differences Backup Diagram



TD estimate: $V_k(s_t) := V_{k-1}(s_t) + \alpha_k \cdot ((r_{t+1} + \gamma V_{k-1}(s_{t+1})) - V_{k-1}(s_t))$

Dynamic Programming Backup Diagram



$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s' \in \mathcal{S}} P(s'|s, a) [R(s, a, s') + \gamma V^\pi(s')]$$

50

Acknowledgments

References

- [1] Kumar Chellapilla, Sidd Puri, and Patrice Simard. High performance convolutional neural networks for document processing. In *Tenth International Workshop on Frontiers in Handwriting Recognition*. Suvisoft, 2006.
- [2] Carl Doersch. Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*, 2016.
- [3] Zhouchen Lin, Risheng Liu, and Zhixun Su. Linearized alternating direction method with adaptive penalty for low-rank representation. In J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24*, pages 612–620. Curran Associates, Inc., 2011.