

$$A = \{a_{0:L-1}\}$$

$$b_k = p(x_k | a_{0:k-1}, z_{1:k})$$

$$p(x_k | x_{k-1}, a_{k-1}); \quad p(z|x)$$

$$r(a,b) = \int_{\mathcal{X}} r(x,a) b(x) dx$$

$$(a) \quad J(b_k, a_{k:k+L-1}) = \mathbb{E}_{z_{k+1:L}} \left\{ \sum_{l=0}^{L-1} r(b_{k+l}, a_{k+l}) + r_T(b_k) \right\}$$

$$(b) \quad J(b_0, a_{0:L-1}) = \mathbb{E}_{z_{1:L}} \left\{ \sum_{l=0}^{L-1} r(b_l, a_l) + r_T(b_0) \right\} =$$

$$\stackrel{\text{linearity}}{=} \mathbb{E}_{z_{1:L}} \left\{ r(b_0, a_0) \right\} + \mathbb{E}_{z_{1:L}} \left\{ \sum_{l=1}^{L-1} r(b_l, a_l) + r_T(b_l) \right\} =$$

$$\stackrel{\text{sufficient statistics}}{=} r(b_0, a_0) + \int_{z_{1:L}} p(z_{1:L} | b_0, a_{0:L-1}) \left(\sum_{l=1}^{L-1} r(b_l, a_l) + r_T(b_l) \right) dz_{1:L} =$$

$$= r(b_0, a) + \int_{z_{1:L}} p(z_1 | b_0, a) p(z_{2:L} | b_0, a_{1:L-1}, z_1) \left(\sum_{l=1}^{L-1} r(b_l, a_l) + \right.$$

chain rule + indep.

$$\left. + r(b_L) \right) dz_{1:L} =$$

$$= r(b_0, a) + \int_{z_1} \int_{z_{2:L}} p(z_1 | b_0, a) p(z_{2:L} | b_0, a_{1:L-1}, z_1) \cdot$$

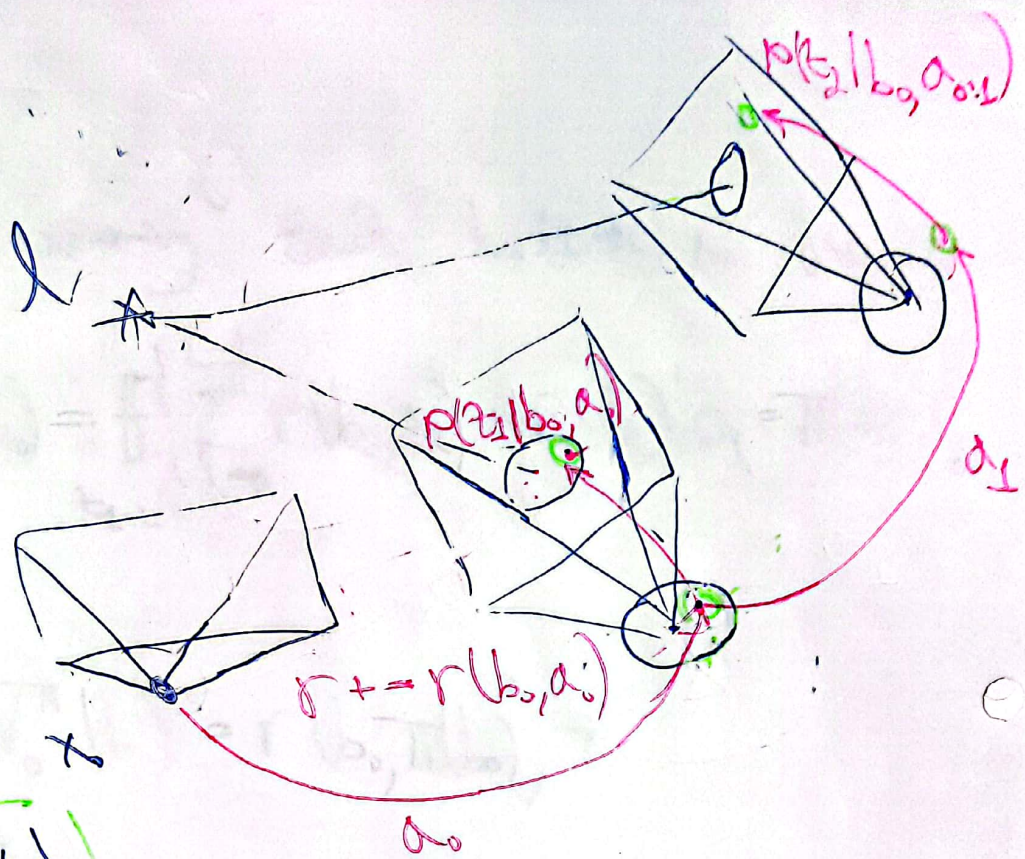
$$\cdot \left(\sum_{l=1}^{L-1} r(b_l, a_l) + r(b_L) \right) dz_{2:L} dz_1 =$$

$$= r(b_0, a) + \int_{z_1} p(z_1 | b_0, a) \int_{z_{2:L}} p(z_{2:L} | b_0, a_{1:L-1}, z_1) \cdot$$

indep.

$$\cdot \left(\sum_{l=1}^{L-1} r(b_l, a_l) + r(b_L) \right) dz_{2:L} dz_1 =$$

$$= r(b_0, a) + \int_{z_1} p(z_1 | b_0, a) T(b | z_1, a_{1:L-1}) dz_1 =$$



$$b_k = \{x_k\}$$

$$J(b_0, a_{0:1}) = r(b_0, a) + \int p(z_1 | b_0, a) \cdot J(b_1, z_1, a_1) dz_1$$

\propto probability
of objective fun

value of
objective fun

$$p(b_1 | z_1) = \underbrace{p(b_0)}_{\text{belief propagation}} p(b_1 | b_0, a_0) p(z_1 | b_1)$$

belief propagation

(c)

(i) assuming, finite horizon; no discount:

$$V_0^\pi(b_0) = \mathbb{E}_{z_{1:L}} \left\{ \sum_{t=0}^{L-1} r(b_t, a_t) + r_T(b_L) \mid a_t = \pi(b_t) \right\}$$

$$(ii) \quad V_0^\pi(b_0) = r(b_0, \pi(b_0)) +$$

$$+ \mathbb{E} \left\{ \sum_{t=1}^{L-1} (r(b_t, a_t) + r_T(b_L)) \mid a_t = \pi(b_t) \right\} =$$

$$= r(b_0, \pi(b_0)) + \int p(z_{1:L} \mid b_0, \pi(b_{0:L-1})) \left[\sum_{t=1}^{L-1} r(b_t, \pi(b_t)) + \right.$$

$$\left. + r(b_L) \right] dz_{1:L} =$$

$$= r(b_0, \pi(b_0)) + \int p(z_1 \mid b_0, \pi(b_0)) p(z_{2:L} \mid b_0, \pi(b_{0:L-1}), z_1) \cdot$$

$$\bullet \left[\sum_{l=1}^{L-1} r(b_l, \pi(b_l)) + r(b_L) \right] dz_{L:L} =$$

$$= r(b_0, \pi(b_0)) + \int_{z_1} p(z_1 | b_0, \pi(b_0)) \int_{z_{2:L}} p(z_{2:L} | b_0, \pi(b_{0:L-1}, z_1)) \cdot$$

$$\bullet \left[\sum_{l=1}^{L-1} r(b_l, \pi(b_l)) + r(b_L) \right] dz_{2:L} dz_1 =$$

$$= r(b_0, \pi(b_0)) + \int_{z_1} p(z_1 | b_0, \pi(b_0)) V_1^\pi(b_1) dz_1$$

(iii) in the objective function formulation J , we assume a set of known actions.

The value function V^π uses a policy $\pi(b)$ to create actions. as such, actions are generated online; and π can be optimized online.