

09/05

$$V = \max_{\pi} \mathbb{E} \sum_{t=0}^{\infty} \gamma^t R(S_t, A_t), \quad \text{if } \gamma < 1, \quad A_t = A(\pi)$$

$$\Rightarrow Q_{k+1}(s, a) = R(s, a) + \gamma \mathbb{E}_{s', a'} V_k(s') \quad s' \sim P(s', a) \rightarrow s' = f(s, a, w)$$

$$V_{k+1}(s) = \max_a Q_{k+1}(s, a) \quad \forall s$$

$$A^*(s) \in \arg \max_a Q(s, a).$$

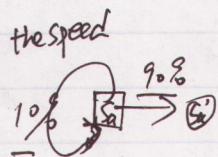
$$\Phi \quad V = \max_{a_1, \dots, a_T} \min_{b_1, \dots, b_T} \sum_{t=1}^T r_t(s_t, a_t, b_t).$$

$$s_{t+1} = f(s_t, a_t, b_t).$$

Reward : $r_t = 1 if caught$

\downarrow iteration counter.

$$EVADER : Q_{k+1}^{evader}(s, a) = 0.9 \eta(s') + 0.1 \eta(s)$$



EVADER :

$$\min_{b_1} \mathbb{E} \sum_{t=1}^{\infty} r_t + \gamma^t$$

PURSUER :

$$\max_{b_1} \mathbb{E} \sum_{t=1}^{\infty} r_t + \gamma^t$$

$$r^a = 90\% \quad \cancel{+ \mathbb{E} V_{k+1}^{evader}(s')}$$

$$V_{k+1}^{evader}(s, a) = \min_a Q_{k+1}(s, a).$$

$\forall s, a$.

$$PURSUER : Q_{k+1}^{purser}(s, b) = 0.6 \eta(s'') + 0.4 \eta(s') + \gamma \mathbb{E} V_{k+1}^{evader}(s'')$$

$$= [0.6 V_{k+1}^{evader}(s'') + 0.4 V_{k+1}^{evader}(s')]$$

$$V_{k+1}^{purser}(s, b) = \max_b Q_{k+1}(s, b)$$

$\forall s, b$

$$s = (x_{evader}, y_{evader}, x_{purser}, y_{purser})$$

V : table of 15^4 , $(15 \times 15)^2$

09/07. Markov Chains.

$X_{t+1} = f(x_t, \pi_{t+1})$.
transaction prob $\pi \geq 0$, $\pi e = 1$.

$$P(X_{t+1} | X_1, \dots, X_t) = P(X_{t+1} | X_t)$$

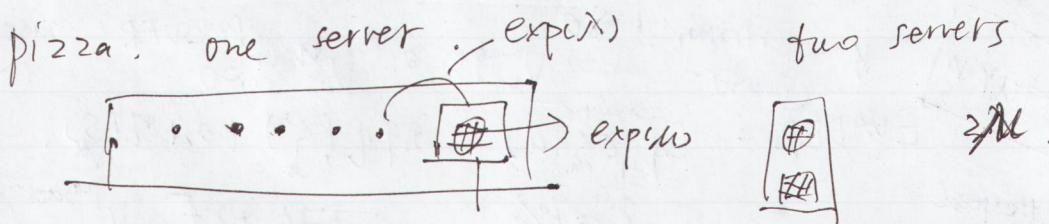
Given $(\pi, p_0, p_1, \dots, p_t, \dots)$

$$\pi_i = \text{prob}(X_0 = i)$$

$$\text{prob}(X_t = j) = \sum_{i=1}^N \underbrace{\text{prob}(X_t = j, X_0 = i)}_{\geq \text{prob}(X_0 = i) \cdot \text{prob}(X_t = j | X_0 = i)} = \sum_{i=1}^N \pi_i p_{ij}$$

$$\Rightarrow P_1 = \pi P_0 \quad \text{Row matrix}$$

$$P_2 = \pi P_0 P_1$$



~~differentiate~~ $\frac{d}{dt} P$

* Mathematical

~~ESC~~ $[\text{ESC}] \alpha [\text{ESC}]$
?2:

expo distribution is memoryless:

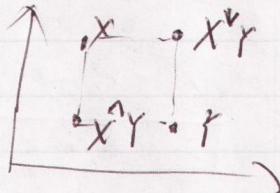
$$P(T > t + \Delta t | T > t) = P(T > \Delta t)$$

info that $T > s$ (we have waiting) is irrelevant
to predict how long to wait

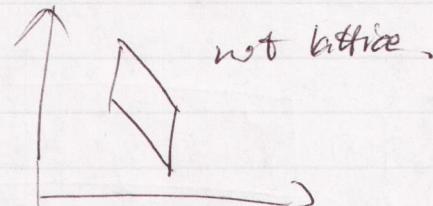
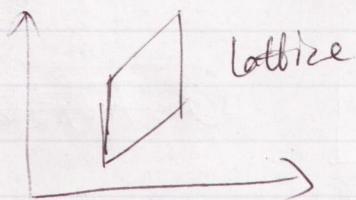
9/19

Lattice
 (\mathbb{R}^2, \leq)

meet $x \wedge y$ $x \wedge y$
join $x \vee y$ $x \vee y$.



remove
↓
...
⇒ still Lattice.



$$(x, y) \in \mathbb{R}^2, x \leq y$$



all format: $x_i - y_i \leq b_i$

$$\text{eg. } L: (x, y, z) : x + y + z \leq 1, x, y, z \geq 0.$$

$$(1, 0, 0) \in L$$

$$(0, 1, 0) \in L$$

$$\text{join } (1, 1, 0) \notin L.$$

$$\left. \begin{array}{l} u = x \\ v = x+y \\ w = x+y+z \end{array} \right\} \rightarrow \left\{ \begin{array}{l} x = u \geq 0 \\ y = v - u \geq 0 \\ z = w - v \geq 0 \\ w \leq 1 \end{array} \right.$$

easy to decide if lattice when 2 variables

g : inventory.

$$g \geq 0.$$

a: sell

$$0 \leq a \leq g.$$

full product x \nearrow product y . planes

$$f(x, y)$$

$$x^+ \geq x^-, y^+ \geq y^-.$$

increasing difference

$$f(x^+, y^-) - f(x^-, y^-) \leq f(x^+, y^+) - f(x^-, y^+)$$

$$f(x^+, y^+) - f(x^+, y^-) \geq f(x^-, y^+) - f(x^-, y^-)$$

$$f(x, y) = x \cdot y$$

Hessian Matrix

$$\geq 0$$

$$\frac{\partial^2 f}{\partial x_i \partial x_j} \geq 0, \forall i, j$$

$$\begin{cases} f(x) \\ g(x) \end{cases}$$

increasing difference

$$\textcircled{2}. \quad h(x) = 0, f(x) + g(x),$$

is convex

$$\frac{\partial^2 h}{\partial x_i \partial x_j} \geq 0$$

① $g(f(x))$ is increasing difference.

supermodular

$$f(x^y) + f(x^y) \geq f(x) + f(y)$$

$x, y \in \text{Lattice}$

$$\text{Proof: } f(x^+, y^+) + f(x^-, y^-) \geq f(x^+, y^-) + f(x^-, y^+)$$

$f_1(x, y), f_2(x, y)$ are supermodular.

$$\text{then } f(x, y) = \theta_1 f_1(x, y) + \theta_2 f_2(x, y), \theta_1, \theta_2 \geq 0$$

is supermodular

in \mathbb{R}^n , supermodular \Leftrightarrow increasing difference.

$$f(\underbrace{u}_{\mathbb{R}^n}, \underbrace{v}_{\mathbb{R}^n}) = u^T v$$

But supermodular \Rightarrow increasing difference ALWAYS

$f(x) = g(a^T x)$ $a \geq 0$, g convex. $f(x)$ is supermodular.

Proof: $\frac{\partial f}{\partial x_i} = g'(a^T x) a_i$

$$\frac{\partial^2 f}{\partial x_i \partial x_j} = \underbrace{g''(a^T x) a_i a_j}_{\geq 0} \geq 0.$$

f is convex if $H \succeq 0$

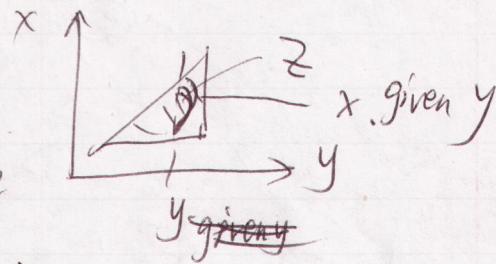
on variable = $f(x)$ must be supermodular

$f(x,y)$ supermodular on $z = \{x,y\}$. z is feasible set

$$g(y) = \sup_{x \in S_y z} f(x,y).$$

assumption (z, \leq) lattice

$\Rightarrow S_z(y)$ lattice



- $g(y)$ supermodularity

$$x^* = \underset{x \in S_y z}{\operatorname{argmax}} f(x,y) \quad (\text{assume unique})$$

* sensitivity to parameters in some models

$y^+ \geq y^- \Rightarrow (x^*)^+ \geq (x^*)^-$ because $f(x,y)$ is lattice

through change parameters.

cannot $x^+ \geq x^-$

$$T(y) = \arg \max_{x \in S_{y,z}} f(x, y)$$

monotonicity y
Isotonicity y

complete monotonic CDF

$$\begin{aligned} E(X) &= \int_0^\infty \underbrace{\Pr[X \geq x]}_{F_X^c(x)} dx = \int_0^\infty 1 - F(x) dx. \\ &= \int_0^\infty 1 dx - \int_0^\infty F(x) dx \end{aligned}$$

$$= \int_0^\infty x f(x) dx = x \Big|_0^\infty - \left(x F(x) \Big|_0^\infty - \int_0^\infty x dF(x) \right).$$

$$= x \Big|_0^\infty - x F(x) \Big|_0^\infty + \int_0^\infty x dF(x)$$

$$= \infty - 0 - \infty + 0$$

$$= \int_0^\infty x f(x) dx.$$

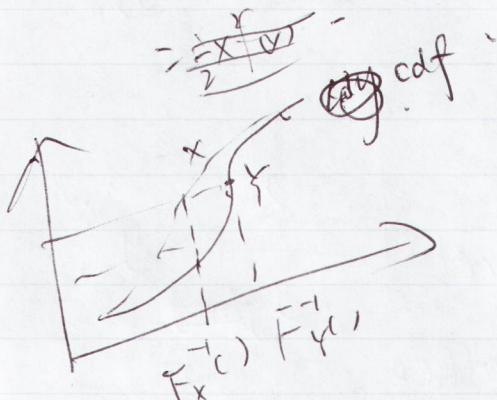
$$X = X^+ - X^- \quad \cancel{E(X) = E(X^+) - E(X^-)}$$

$$E[X^+] = \int_0^\infty x^+ f(x) dx = \int_0^\infty x f(x) dx$$

$$E[X^-] = \int_0^\infty x^- f(x) dx = \int_0^\infty (-x) f(x) dx$$

$$X^+ = \begin{cases} X, & x \geq 0 \\ 0, & x < 0 \end{cases}, \quad X^- = \begin{cases} -x, & x \leq 0 \\ 0, & x > 0 \end{cases}$$

$$E(X) = E[X^+] - E[X^-]$$



$$X \leq_{st} Y \Rightarrow E[X] \leq E[Y]$$

$X \leq_{st} Y \Rightarrow h(x) \leq_{st} h(y)$. h : increasing function.

$$\Rightarrow E[h(X)] \leq E[h(Y)]$$

$$X'_\alpha = \underset{\text{current}}{\text{stochastic}}$$

$$X'_{S_1} \leq_{st} X'_{S_2}$$

$$E[V(X')] \leq_{x' \mid S_2} E[V(X)]$$

$$S_1 \leq S_2$$

Stochastic: First Order Dominance:

Isotone

$$\begin{aligned} X &\leq_{st} Y \\ \xrightarrow{\quad} F_X^C(t) &= 1 - \text{Prob}(X \leq t) = \text{Prob}(X > t) \\ F_X^C(t) &\leq F_Y^C(t). \end{aligned}$$

"stochastic greater or less than".

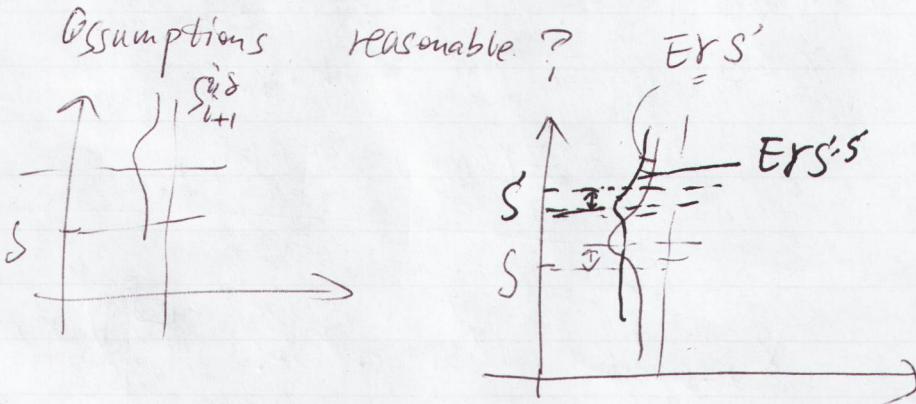
For $\{X_\alpha\}$, X_α is stochastic increasing in α .

$$X_{\alpha_1} \leq X_{\alpha_2}, \forall \alpha_1 \leq \alpha_2$$

- Exp Random Value are sto decreasing

ISE 416.

~~Sep 28~~, Oct 3.



$$S_{t+1} - S_t = \underset{\text{Not}}{(M_t - S_t) \alpha t + \sigma(t, S_t) \cdot \sqrt{t}} N(0, 1)$$

continuous \rightarrow discrete : $dS_t \rightarrow S_{t+1} - S_t$
 $dt \rightarrow \Delta t$
 $d\beta_t \rightarrow \sqrt{\Delta t} N(0, 1)$

ISE 41B

DP: Infinite horizon: $\max \sum_{t=0}^{\infty} r^t R_t$. $R_t = R_t(s_t, a_t)$

$$V(s) = \max_a Q(s, a).$$

$$Q(s, a) = R(s, a) + r \mathbb{E} V(s')$$

$$V(s) = \max_a \left[R(s, a) + r \mathbb{E}_{s' | s, a} V(s') \right]$$

Algorithm: Value iteration: (K is. ite 次數)

Set V_0 ~~arbitrarily~~. arbitrarily.

$$V_{k+1}(s) = \max_a [R(s, a) + r \mathbb{E}_{s' | s, a} V_k(s')]$$

$$\text{hope: } V_{k+1} \rightarrow V^\infty$$

Expect value of $\tilde{\pi}$:

$$V^{\tilde{\pi}}(s_0) = V(\tilde{\pi})|_{s_0} = \mathbb{E} \left[\sum_{t=0}^{\infty} r^t R_t(s_0, A_t^{\tilde{\pi}}(s_t)) | s_0 \right]$$

$$V^{\tilde{\pi}}(s) = R(s, A^{\tilde{\pi}}(s)) + r \mathbb{E} V^{\tilde{\pi}}(s')$$

$$V_{k+1}^{\tilde{\pi}}(s) = R(s, A^{\tilde{\pi}}(s)) + r \mathbb{E}_{\substack{s' | s, a^{\tilde{\pi}}(s) \\ s' | s, a = A^{\tilde{\pi}}(s)}} V^{\tilde{\pi}}(s')$$

$$\boxed{V_{k+1} = T V_k} \quad T: \text{Bellman operation}.$$

Value iteration algorithm

$$V_{k+1} = T^{\tilde{\pi}} V_k$$

↑ operator ↑ function

$(T^{\tilde{\pi}}$ is no just a number)

want $\rightarrow V = TV$. fixed-point solution.

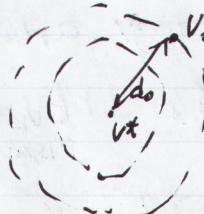
$$V_{k+1} = T^{(k+1)} V_0.$$

↓ function ↑ super norm
 $\|V\|_{\sup} = \sup_S (|V(s)|).$

$$\|V_2 - V_1\| = \sup_S |V_2(s) - V_1(s)|.$$

$$\|V_{k+1} - V_k\| \rightarrow 0.$$

$$d_0 = \|V^* - V_0\|$$



$$\|V_{k+1} - V^*\| < \|V_k - V^*\|$$

$$V_{k+1} = TV_k.$$

T : contractive.

T : Non expansive. $\|V_{k+1} - V^*\| \leq \|V_k - V^*\|$ may not be converged.

proof $\|TV_k - TV^*\| \leq r \|V_k - V^*\|$

\Downarrow

$$\|V_{k+1} - V^*\| \leq r \|V_k - V^*\|.$$

$$\leq r^{k+1} \|V_0 - V^*\|$$

$$d := \|J - J'\|.$$

$$J - d \leq J' \leq J + d.$$

By monotonicity of T , $TJ - rd \leq TJ' \leq TJ + rd$.

~~$$-rd \leq TJ' - TJ \leq rd$$~~

$$\|TJ' - TJ\| \leq rd = r\|J - J'\|.$$

ZSE 416 .
1%

$$TJ^* = J^*, \\ \|TJ - TJ'\| \leq \delta \|J - J'\|.$$

$TJ - J \equiv$ Bellman Residual .

converge .

$$\|J^* - J\| \leq \frac{1}{1-r} \|TJ - J\|. \quad \text{threshold}$$

$\uparrow \quad \downarrow$
 $J^* = TJ^*$ can be computed.

$$\|J^* - TJ\| = \|TJ^* - TJ\| \leq r \|J^* - J\| \leq \frac{r}{1-r} \|TJ - J\|. \quad \text{tidy diff.}$$

$$\tilde{\pi} \rightarrow \tilde{J}_{\tilde{\pi}} ? \quad T^{\tilde{\pi}} \tilde{J}_{\tilde{\pi}} - \tilde{J}_{\tilde{\pi}} \quad \delta = \|T^{\tilde{\pi}} \tilde{J}_{\tilde{\pi}} - \tilde{J}_{\tilde{\pi}}\|$$

$$\|\tilde{J}_{\tilde{\pi}} - J_{\tilde{\pi}}^*\| \leq \frac{1}{1-r} \|T^{\tilde{\pi}} \tilde{J}_{\tilde{\pi}} - \tilde{J}_{\tilde{\pi}}\|.$$

$$\leftarrow \overset{*}{\rightarrow} | \quad \leftarrow \rightarrow \quad \begin{matrix} \uparrow \\ \tilde{J}_{\tilde{\pi}_1} \end{matrix} \quad \begin{matrix} \uparrow \\ \tilde{J}_{\tilde{\pi}_2} \end{matrix}$$

$$\tilde{J} \rightarrow \pi_{\tilde{J}} \rightarrow J_{\pi_{\tilde{J}}} \rightarrow \tilde{J}_{\pi_{\tilde{J}}}. \quad \begin{matrix} \uparrow \\ \text{can't again.} \end{matrix}$$

$$\hookrightarrow = \sup_{\pi} \mathbb{E} \sum_{t=0}^{\infty} r^t R(s, \pi(s)).$$

$J = \pi J$. can't solve this exact value .

ISE 416. 10/12.

Policy Iteration:

* Exact solve the equation as written

never worse than value iteration

* Approximate

not necessarily true. $\xrightarrow{\text{way}} \text{"Bellman" method}$

Iteration k. 1. Policy evaluation : $J_{\pi_k} = T^{\pi_k} J_{\pi_k}$.

2. Policy improvement : $s : \max_a R(s, a) + \gamma \mathbb{E}_{s' \sim P(s)} [J_{\pi_k}(s')]$

$$\begin{aligned} J' &\geq J \\ \Rightarrow TJ' &\geq TJ \quad \text{monotonicity} \end{aligned} \quad \xrightarrow{T\pi_{k+1} J_{\pi_k} = T J_{\pi_k}} (T J_{\pi_k} \neq T_{\pi_k} J_{\pi_k}).$$

$$\pi_k(s) = \arg \max_a \mathbb{E}_{s'|s,a} [R(s,a) + \gamma J_{\pi_k}(s')]$$

$$J_{\pi_0}(s) = \mathbb{E} \sum_{t=0}^{\infty} r^t R(s, \pi_0(s)) = \mathbb{E}_a [R(s, \pi_0(s))] + \gamma \mathbb{E}_{s'|s,a} J_{\pi_0}(s')$$

$$\begin{aligned} s=1 &: V_{\pi_0} = \mathbb{E}_{\pi_0} + r \mathbb{E}_{\pi_0} [P_{\pi_0} V_{\pi_0}] \\ s=N & \end{aligned}$$

$$P_{\pi_0} = \begin{array}{c|c} & \text{row of } \pi_0(i) \\ \hline \end{array}$$

$$(I - rP_{\pi_0}) X = \mathbb{E}_{\pi_0}$$

$$\Rightarrow X = (I - rP_{\pi_0})^{-1} \mathbb{E}_{\pi_0}$$

$$3. T_{K+1} \equiv T_K, \Rightarrow T_{\pi_K} J_{\pi_K} = T_{\pi_{K+1}} J_{\pi_K} = J_{\pi_K}.$$

3.8 Approximate Policy Iteration

J_K $\xrightarrow{\text{real value}}$ because of inverse of V matrix
 $|J_K - J_{\pi^K}| \leq \delta \Rightarrow \|T_{\pi_K} J_K - T J_K\| \leq \epsilon$.
 ↗ approximate value of T_K .

$$J_K = (1-\lambda) \tilde{J}_{\pi_K} + \lambda \tilde{J}_{\pi_{K+1}}$$

10/189. Value Iter : find V : $V = TV \Rightarrow T_{\pi^V} = TV_K$ until convergence.

Policy Iter : find π : evaluate: $V^\pi = T^\pi V^\pi$.

$$V_{k+1}^\pi = T^\pi V_k^\pi \quad A_t = A^\pi(s_t) \text{ improve: find better } \pi^{t+1}$$

Linear Optimization :
 • finit MDP, you get exact V .
 however, you solve only small problem.

• Easy to add constraint.

finit MDP. on s, a .

$$\text{find policy: } \pi = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$$

$$A_t = A^\pi(s_t) = P_{S_t j} = \text{Prob}(S_{t+1} = j \mid S_t = s) \quad A^\pi(s)$$

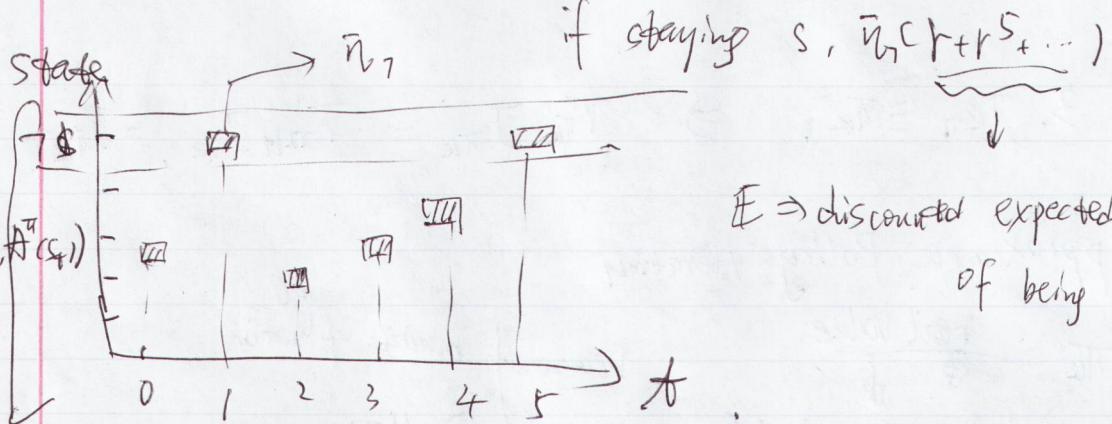
$$V^\pi = \mathbb{E} \sum_{t=0}^{\infty} r^t R(S_t, A^\pi(s_t))$$

$$\tilde{r}_s^\pi = \mathbb{E} R(S, A^\pi(s), S_{t+1})$$

$$V^\pi = (I - \gamma P_\pi^\top)^{-1} \pi_{\pi_0}$$

$$= \sum_{j=1}^{\infty} P_{S_j}^\pi \cdot R(S, A^\pi(s), j)$$

$$\sum_{t=0}^{\infty} \gamma^t = \frac{1}{1-\gamma} = (1-\rho)^{-1}$$



$\mathbb{E} \Rightarrow$ discounted expected count
of being in state j

$$V_s^\pi = \sum_{j=1}^m \bar{n}_j^\pi \cdot \mathbb{E} \left\{ \begin{array}{l} \text{discounted} \\ \text{count of} \\ \text{being at} \\ \text{state } j \end{array} \right\}$$

given starting state s and policy π $(\sum_m = 1)$
 $(\frac{1}{1+\gamma})$ [discounted fraction of time on a state j]

$$V^\pi(s) = \text{row } s \text{ of } (I - \gamma P^\pi)^{-1} \bar{n}^\pi = \frac{M}{1-\gamma} \bar{n}^\pi$$

$$M = \frac{(I - \gamma P^\pi)^\gamma}{(\frac{1}{1-\gamma})} = (1-\gamma) (I - \gamma P^\pi)^{-1}$$

normalize rows of M . = discount fraction of time in state 1, 2, ..., m .

$$M_j \quad \begin{matrix} 1 & 2 & \dots & m \\ \hline 1 & \dots & j & \dots & m \end{matrix}$$

$$S \quad \begin{matrix} M \\ \hline M \end{matrix} \quad m \times m$$

M : OCCUPATION measure.

like probability but not exactly real.

$$I + \gamma P + \gamma^2 P^2 + \gamma^3 P^3 + \dots = (I - \gamma P)^{-1}$$

10/31.

$$Q_{k+1}(s, a) = \gamma r(s, a) + \max_{s' \in S} V_k(s')$$
$$= n(s, a) + \max_{a' \in A(s)} Q_k(s', a').$$

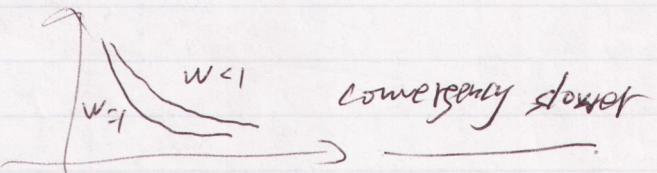
$$Q_{k+1}(s, a) = (1 - \alpha(s, a)) Q_k(s, a).$$

$$+ \alpha(s, a) [n(s, a) + \max_{a' \in A(s)} Q_k(s', a')].$$

$\alpha(s, a) \sim N(s, a)$: # times visiting (s, a) .

$$\alpha(s, a) = \left[\frac{1}{N(s, a)} \right]^w$$
, polynomial step size

$$\begin{cases} \sum_k \alpha_k^2 < \infty, & w \in [0, 1], \\ \sum_k \alpha_k = \infty. & \end{cases}$$



$$(s_0, a_0) \xrightarrow{\text{code}} (s_1, a_1) \xrightarrow{\text{code}} (s_2, a_2).$$

visit all state-action pair
as often.

with probability $1 - \epsilon$: $a_k = \arg \max Q_k(s_k, a_k)$

ϵ : $a_k \in A(s_k)$ randomly

④ If the step size is too big

convergence at a bad point

check shaking, window(max-min) ≤ threshold

0.8

Asynchronous Q-learning

Double Q-learning

Speedy Q-learning

1/2. feature.

$$V(s) = \varphi_1(s) \cdot \theta_1 + \dots + \varphi_m(s) \cdot \theta_m$$

finite MDP, 5 states.

$$\begin{aligned} \varphi_1(s) &= I(s=1) \\ &\vdots \\ \varphi_5(s) &= I(s=5) \end{aligned} \quad \left\{ \begin{array}{l} m \ll |S| \\ \end{array} \right.$$

basic function = $\{\varphi_k\}_{1 \leq k \leq m}$.

$$f(x) \approx \sum_{k=0}^m w_k \sin(kx) + w_k^{(2)} \cos(kx)$$

~~w₀⁽²⁾~~

$\simeq \varphi_k$

$$V_{\varphi, \theta}(s) = \varphi_1(s) \cdot \theta_1 + \dots + \varphi_m(s) \cdot \theta_m$$

$$\left\{ \begin{matrix} s_i \\ x \end{matrix}, \begin{matrix} V(s_i) \\ y \end{matrix} \right\}_{i=1 \dots n} \Rightarrow \theta: \min \frac{1}{2} \| \varepsilon \|_2^2 = \frac{1}{2} \| V^{\text{target}} - \Phi \theta \|_2^2$$

$$\varepsilon_i = V(s_i) - V_{\varphi, \theta}(s_i)$$

$$\Phi = \begin{bmatrix} \varphi_1(s_1) & \dots & \varphi_m(s_1) \\ \vdots & \ddots & \vdots \\ \varphi_1(s_N) & \dots & \varphi_m(s_N) \end{bmatrix}$$

$$\theta = (\Phi^\top \Phi)^{-1} (\Phi^\top V^{\text{target}})$$

$$\min \| \varepsilon \|_2^2. \quad \| V \|_W = \sqrt{V^\top W V}$$

$\underbrace{\text{dim } N}_{\text{dim } V}$

diag $\begin{bmatrix} w_1 & w_2 & \dots & w_N \end{bmatrix}$ can be determined by
the occupation percentage prob.

1/5

$$V^{\pi} = T^{\pi} V^{\pi}$$

$$V_{(x)}^{\pi} = \mathbb{E}_{y|x, \beta^{\pi}(x)} (r + \gamma V^{\pi}(y))$$

$$V_{(x)}^{\pi} = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_{\pi}(x_t, \beta^{\pi}(x_t, x_{t+1})) \mid x_0 = x \right]$$

$$V^{\pi} = T^{\pi} V^{\pi} \quad V^{\pi}(s) = \bar{\Phi}(s)\theta = \phi_1(s)\theta_1 + \dots + \phi_m(s)\theta_m$$

$$\min_{\theta} \|\bar{\Phi}\theta - T^{\pi} \bar{\Phi}\theta\|_2 \quad \leftarrow \min_{\theta} \|V^{\pi} - T^{\pi} V^{\pi}\|_{\text{sup}}$$
$$\bar{n} + \gamma P^{\pi} \bar{\Phi} \theta.$$

$$V(s_1) \sim \bar{\Phi}(s_1) + P^{\pi} \bar{\Phi} \theta.$$

$$\min_{\theta} \|\underbrace{(I - \gamma P^{\pi}) \bar{\Phi}}_{\Psi} \theta - \bar{n}\|_2 \quad \text{estimate } P^{\pi}.$$

$$(\bar{\Phi} \bar{\Phi}^T)^{-1} (\bar{\Phi}^T \bar{n})$$

$$\Rightarrow \|\bar{\Phi}\theta - \Pi_{\bar{\Phi}}(T^{\pi} \bar{\Phi}\theta)\| \quad \begin{matrix} \nearrow \\ \cancel{\text{Bellman}} \end{matrix} \quad \begin{matrix} \nearrow \\ \text{Bellman Residual minimization} \end{matrix}$$

1/9.

V_K : approximation $\Leftrightarrow \theta_K$.

$$\text{By } V_K(y) = \varphi_1(y)\theta_{k1} + \dots + \varphi_m(y)\theta_{km}.$$

Sample $x_i \sim \mu, i=1, \dots, N$.

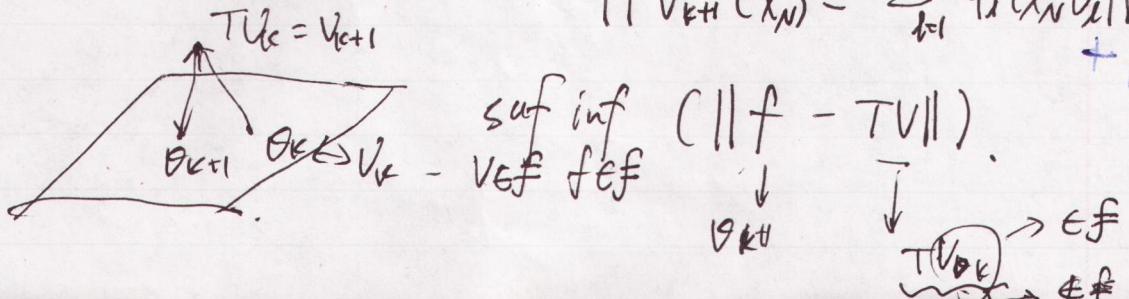
$$x_i \rightarrow \tilde{V}_{k+1}(x_i) = \max_a \mathbb{E}_{y|x_i} [\eta(x_i, a, y) + \delta V_k(y)]$$

$$\vdots \quad \tilde{V}_{k+1}^M(x_i) = \max_a \frac{1}{M} \sum_{j=1}^M [\eta(x_i, a, y_{ij}) + \delta V_k(y_{ij})]$$

$$x_N \rightarrow \tilde{V}_{k+1}^M(x_N) = \max_a \frac{1}{M} \sum_{j=1}^M [\eta(x_N, a, y_{Nj}) + \delta V_k(y_{Nj})].$$

$$[(x_1, \tilde{V}_{k+1}^M(x_1)), \dots, (x_N, \tilde{V}_{k+1}^M(x_N))]$$

$$\min_{\theta_{k+1}} \|\tilde{V}_{k+1}^M - \tilde{\theta}\| = \left\| \begin{array}{l} \tilde{V}_{k+1}(x_1) - \sum_{i=1}^m \varphi_i(x_1) \theta_i \\ \vdots \\ \tilde{V}_{k+1}(x_N) - \sum_{i=1}^m \varphi_i(x_N) \theta_i \end{array} \right\| + \|\theta\| \leq A$$



$$\textcircled{D} \quad V_S = \max_a \pi_{S,a} + \gamma \sum_{s'} P_{S,s'}^a V$$

Stochastic policy

$$y_{S,a} = \text{Prob}(a_t = a | S_t = s)$$

$$= \max_a \pi_{S,a} + \gamma \sum_j P_{S,j}^a V_j$$

$$\left\{ \begin{array}{l} V_S \geq \pi_{S,a} + \gamma \sum_j P_{S,j}^a V_j \quad \forall a, s \\ \text{aux constraint} \end{array} \right.$$

$$\text{objective: } \min_S \sum_{S \geq 0} w_S V_S$$

dual value $\neq 0$, the optimal a at states.

for each state
sum of dual value is

$$M(S).$$

$$\max_x \min_{R' \in Q} \sum_{S,a} R'_{S,a} x_{S,a}$$

$$\text{s.t. } \sum_a x_{S,a} - \gamma \sum_i \sum_a P_{i,S}^a x_{i,a} = \bar{w}_S \quad \forall S.$$

$$x_{S,a} \geq 0.$$

Robust MDP.

$$\mathcal{L} = \mathbb{E}(\bar{R}, \Sigma). \\ (\cdot, \bar{R}, \Sigma) \quad \text{s.t. } (y - \bar{R})^T \Sigma^{-1} (y - \bar{R}) \leq \alpha^2.$$

$$\max_{\pi} \min_{R'} \mathbb{E} \sum_{t=0}^{\infty} \gamma^t R'(S_t, a_t).$$

$$\min_{y \in Q} x^T y = \min_{y \in Q} x^T y$$

$$\text{s.t. } (y - \bar{R})^T \Sigma^{-1} (y - \bar{R}) \leq \alpha^2.$$

$$\| \Sigma^{1/2} (y - \bar{R}) \|_2 \leq \alpha.$$

$$\max_{\pi} -\alpha \sqrt{x^T \Sigma x} + \bar{R}^T x$$

$$\max_{\pi} \sum_{S,a} \bar{P}_{S,a} x_{S,a} - \alpha$$

$$\rightarrow \alpha \sqrt{\sum_{S,a} \sum_{S',a'} x_{S,a} \sum_{S',a'} \bar{P}_{S,a}^a x_{S',a'}}$$

$$\text{s.t. } \sum_a x_{S,a} - \gamma \sum_i \sum_a P_{i,S}^a x_{i,a} = \bar{w}_S$$

$$x_{S,a} \geq 0.$$

10/20

CODE

$$y_{t+1} = f(s_t, a_t)$$
$$\hat{p}_{t+1} = n(s_t, a_t).$$

$(s_t, a_t, s_{t+1}, a_{t+1})$

MODEL - BASED.

MDP: R, P . n : # of sampling

$$P_{sj}^a = \frac{1}{M} \sum_{m=1}^M \mathbb{I}\{s' = j | s, a\} = \frac{N_{saj}}{M} \xrightarrow{M \rightarrow \infty} P(j | s, a).$$

:

$$\bar{R}_{sa} = \frac{1}{M} \sum_{m=1}^M \bar{R}(s, a) \xrightarrow{M \rightarrow \infty} \bar{R}(s, a)$$

$\rightarrow \frac{1}{n}$ optimal for \hat{P}_m, \hat{R}_m

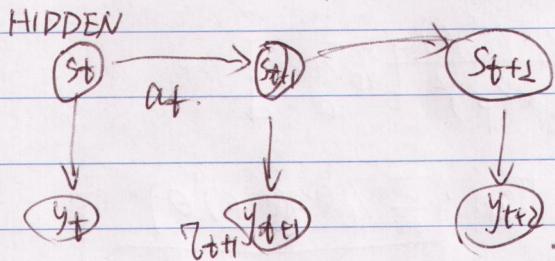
$$\rightarrow V^{\frac{1}{n}}$$

11/21.

POMDP : Partially observable MDP.

MDP : $x_t, a_t \rightarrow x_{t+1}, a_{t+1}$.

POMDP :



$$P(S_t = i | y_0, y_1, \dots, y_t).$$

$$P(S_{t+1} = i | y_0, y_1, \dots, y_{t+1}) = f(P(S_t = i | y_0, \dots, y_t), y_{t+1}).$$

$$P^a, R^a$$

$$Y^a_s = \text{Prob}(y_t = k | S_t = s, A_t = a).$$

$$\max_{\pi \in \mathcal{F}} \mathbb{E} \sum_{t=0}^{\infty} R_t.$$

$$\max_{\theta \in \Theta} \mathbb{E} \sum_{t=0}^{\infty} R_t. \rightarrow \max_{\theta \in \Theta} J^{\pi^{\theta}}(s_0) = f(\theta)$$

Gitts policy

$$\begin{bmatrix} \phi_1(s, a) \\ \phi_2(s, a) \\ \vdots \\ \phi_m(s, a) \end{bmatrix}$$

$$Q(s, a) \cong \phi(s, a)^T \theta \quad \theta \in \mathbb{R}$$

$$\text{Prob}(a|s) = \pi(s, a) = \frac{\exp(\phi(s, a)^T \theta)}{\sum_a \exp(\phi(s, a)^T \theta)}$$

$$[\pi(s, 1), \dots, \pi(s, m)] \quad a \sim \pi,$$

$$\begin{aligned}
 \nabla_{\theta} \pi(s, a) &= \nabla_{\theta} \left[\frac{1}{\sum_{a'} \exp(\phi(s, a')^T \theta)} \right] \\
 &= \left[\nabla_{\theta} \frac{1}{\sum_{a'} \exp} \right] \exp(\phi(s, a)^T \theta) + \frac{1}{\sum_{a'} \exp} \nabla_{\theta} \exp(\phi(s, a)^T \theta) \\
 &= \frac{-\exp(\phi(s, a)^T \theta)}{\left(\sum_{a'} \exp \right)^2} \left[\nabla_{\theta} \sum_{a'} \exp \right] + \dots \\
 &= -\pi(s, a) \cdot \left[\frac{\sum_{a'} \nabla_{\theta} \exp(\phi(s, a')^T \theta)}{\sum_{a''} \exp(s, a'')} \right] + \frac{1}{\sum_{a'} \exp} \left[\phi(s, a) \exp(\phi(s, a)^T \theta) \right] \\
 &= -\pi(s, a) \left[\frac{\sum_{a'} \phi(s, a') \exp(\phi(s, a')^T \theta)}{\sum_{a''} \exp(\phi(s, a')^T \theta)} \right] + \phi(s, a) \cdot \pi(s, a) \\
 &= \pi(s, a) \left[\phi(s, a) - \sum_{a'} \phi(s, a') \pi(s, a') \right]
 \end{aligned}$$

$$\nabla_{\theta} \log \pi(s, a) = \frac{1}{\pi(s, a)} \nabla_{\theta} (\pi(s, a)) = [\phi(s, a) - \sum_{a'} \phi(s, a') \pi(s, a')]$$

$\nabla_{\theta} E_{\theta} f(x) \quad x \sim P_{\theta}$.

$$\begin{aligned}
 &\nabla_{\theta} \int f(x) p(x|\theta) dx \\
 &= \int \nabla_{\theta} [f(x) p(x|\theta)] dx. \\
 &= \int f(x) \nabla_{\theta} (p(x|\theta)) dx. \\
 &= \int f(x) p(x|\theta) \nabla_{\theta} \log p(x|\theta) dx \\
 &= E_{\theta} [f(x) \nabla_{\theta} \log p(x|\theta)],
 \end{aligned}$$

SCORE

METHOD.

policy search

$$\sim \frac{1}{N} \sum_{m=1}^N [f(x_m) \nabla_{\theta} \log p(x_m|\theta)]$$

~~fisher matrix
(information)~~

Homework 3

Fuxin Xu.

Question 1. Show that a decreasing convex function of an increasing submodular function is supermodular. For simplicity, you can assume the functions are twice differentiable.

Answer. Set: $g(x)$ is a decreasing convex function.

$f(x)$ is an increasing submodular function.

Because $g(x)$ is decreasing, $f(x)$ is increasing.

Thus $g(f(x))$ is decreasing, $g'(f(x)) \leq 0$.

Because $f(x)$ is submodular,

thus $\frac{\partial^2 f(x)}{\partial x_i \partial x_j} \leq 0, \forall i \neq j$.

Because $f(x)$ is increasing thus $f(x) \geq 0$

Because $g(x)$ is convex. Thus $g''(f(x)) \geq 0$.

$$\frac{\partial f(x)}{\partial x_i}, \frac{\partial f(x)}{\partial x_j}$$

$$\text{Thus: } \frac{\partial^2 g(f(x))}{\partial x_i \partial x_j} = \frac{\partial f(x) \cdot g'(f(x))}{\partial x_j} = \underbrace{\frac{\partial f(x)}{\partial x_i \partial x_j}}_{\leq 0} \cdot \underbrace{g'(f(x))}_{\leq 0} + \underbrace{f'(x) \cdot f''(x)}_{\geq 0} \cdot \underbrace{g''(f(x))}_{\geq 0} \geq 0$$

$$\geq 0. \quad \forall i \neq j.$$

Thus, it's a supermodular.

Question 2. A family of random variables $\{X_s\}$ parameterized by s is said to be stochastically convex in s if their complementary cumulative distribution function $F_s^c(t) = \text{Prob}(X_s > t)$ is convex in the parameter s for each fixed t . Show that the expectation function $h(s) = \mathbb{E}[X_s]$ is convex in s .

Answer.

$$h(s) = \mathbb{E}[X_s] = \int_0^\infty F_s^c(x) dx$$

$$\begin{aligned} \frac{\partial^2 h(s)}{\partial s^2} &= \frac{\partial^2 \int_0^\infty F_s^c(x) dx}{\partial s^2} = \frac{\partial \int_0^\infty \frac{\partial F_s^c(x)}{\partial s} dx}{\partial s} \\ &= \int_0^\infty \frac{\partial^2 F_s^c(x)}{\partial s^2} dx \end{aligned}$$

Because $F_s^c(x)$ is convex, thus $\frac{\partial^2 F_s^c(x)}{\partial s^2} \geq 0$.

Thus, $\frac{\partial^2 h(s)}{\partial s^2} \geq 0$, $h(s) = \mathbb{E}[X_s]$ is convex in s .

Project

$$\begin{aligned} \min \quad & \frac{1}{2} w^T w + C \sum_{i=1}^N [\tau_m \xi_i^+ + (1-\tau_m) \xi_i^-] \\ \text{s.t.} \quad & y - xw \leq \varepsilon + \xi_i^+ \\ & -y + xw \leq \varepsilon + \xi_i^- \\ & \xi_i^+ \geq 0, \quad \xi_i^- \geq 0 \end{aligned}$$

$y - xw \leq \varepsilon + \xi_i^+ \quad \Rightarrow \quad xw - y + \varepsilon + \xi_i^+ \geq 0$

$-y + xw \leq \varepsilon + \xi_i^- \quad \Rightarrow \quad -xw + y + \varepsilon + \xi_i^- \geq 0$

Lagrangian Duality:

$$\begin{aligned} L(w, \xi^+, \xi^-, \alpha^+, \alpha^-, \beta^+, \beta^-) = & \frac{1}{2} w^T w + C \sum_{i=1}^N [\tau_m \xi_i^+ + (1-\tau_m) \xi_i^-] \\ & - \alpha^{+T} (xw - y + \varepsilon + \xi_i^+) - \alpha^{-T} (-xw + y + \varepsilon + \xi_i^-) - \beta^{+T} \xi_i^+ - \beta^{-T} \xi_i^- \end{aligned}$$

$$\frac{\partial L}{\partial w} = w - X^T \alpha^+ + X^T \alpha^- \stackrel{\text{set}}{=} 0 \Rightarrow w = X^T \alpha^+ - X^T \alpha^-$$

$$\frac{\partial L}{\partial \xi^+} = C \cdot \tau_m \cdot I - \alpha^+ - \beta^+, \quad \frac{\partial L}{\partial \xi^-} = C \cdot (1-\tau_m) - \alpha^- - \beta^-.$$

$$\begin{aligned} L = & \frac{1}{2} (X^T \alpha^+ - X^T \alpha^-)^T (X^T \alpha^+ - X^T \alpha^-) + C \sum_{i=1}^N [\tau_m \xi_i^+ + (1-\tau_m) \xi_i^-] \\ & - \alpha^{+T} [X(X^T \alpha^+ - X^T \alpha^-)] - \alpha^{-T} [-X(X^T \alpha^+ - X^T \alpha^-)] - \beta^{+T} (-y + \varepsilon + \xi_i^+) \\ & - \beta^{-T} (y + \varepsilon + \xi_i^-) - \beta^{+T} \xi_i^+ - \beta^{-T} \xi_i^- \end{aligned}$$

$$\begin{aligned} = & \frac{1}{2} (X^T \alpha^+ - X^T \alpha^-)^T (X^T \alpha^+ - X^T \alpha^-) + C \sum_{i=1}^N [\tau_m \xi_i^+ + (1-\tau_m) \xi_i^-] \\ & - (\alpha^{+T} X - \alpha^{-T} X) (X^T \alpha^+ - X^T \alpha^-) - \alpha^{+T} (-y + \varepsilon + \xi_i^+) - \alpha^{-T} (y + \varepsilon + \xi_i^-) \\ & - \beta^{+T} \xi_i^+ - \beta^{-T} \xi_i^- \end{aligned}$$

$$\begin{aligned} = & -\frac{1}{2} (X^T \alpha^+ - X^T \alpha^-)^T (X^T \alpha^+ - X^T \alpha^-) + C \sum_{i=1}^N [\tau_m \xi_i^+ + (1-\tau_m) \xi_i^-] \\ & - (X^T \alpha^+ - X^T \alpha^-)^T (X^T \alpha^+ - X^T \alpha^-) - \alpha^{+T} (-y + \varepsilon + \xi_i^+) - \alpha^{-T} (y + \varepsilon + \xi_i^-) \\ & - \beta^{+T} \xi_i^+ - \beta^{-T} \xi_i^- \end{aligned}$$

$$\begin{aligned} = & -\frac{1}{2} (X^T \alpha^+ - X^T \alpha^-)^T (X^T \alpha^+ - X^T \alpha^-) + C \cdot \tau_m \cdot I^T \cdot \xi_i^+ \\ & + C \cdot (1-\tau_m) \cdot I^T \cdot \xi_i^- + C \cdot (1-\tau_m) \cdot I^T \cdot \xi_i^- - \beta^{+T} \xi_i^+ - \beta^{-T} \xi_i^- \end{aligned}$$

$$\begin{aligned} = & -\frac{1}{2} (X^T \alpha^+ - X^T \alpha^-)^T (X^T \alpha^+ - X^T \alpha^-) + [C \cdot \tau_m \cdot I^T - \alpha^{+T} - \beta^{+T}] \xi_i^+ \\ & + [C \cdot (1-\tau_m) \cdot I^T - \alpha^{-T} - \beta^{-T}] \xi_i^- + \alpha^{+T} (y - \varepsilon) - \alpha^{-T} (y + \varepsilon) \end{aligned}$$

$$\begin{aligned} = & -\frac{1}{2} \alpha^{+T} X X^T \alpha^- - \frac{1}{2} \alpha^{-T} X X^T \alpha^+ + \alpha^{+T} X X^T \alpha^- + \alpha^{+T} (y - \varepsilon) - \alpha^{-T} (y + \varepsilon) \end{aligned}$$

$$\frac{1}{2} \alpha^{+T} K \alpha^- + \frac{1}{2} \alpha^{-T} K \alpha^+$$

$$K = X X^T$$

$$\max \quad \underbrace{\alpha^{+T} X X^T \alpha^-}_{\text{K}} - \frac{1}{2} \alpha^{+T} X X^T \alpha^+ - \frac{1}{2} \alpha^{-T} X X^T \alpha^- + \alpha^{+T} (y - \varepsilon) - \alpha^{-T} (y + \varepsilon)$$

s.t. $\alpha_i^+ \in [0, C \tau_m], \forall i$

$$\alpha_i^- \in [0, C(1 - \tau_m)], \forall i$$

$$\cancel{\frac{\partial L}{\partial \alpha^+}} = \cancel{\alpha^{+T} \alpha^-} - \cancel{\alpha^{+T} \alpha^+} + y - \varepsilon.$$

$$\cancel{\frac{\partial L}{\partial \alpha^-}} = \cancel{\alpha^{-T} \alpha^+} - \cancel{\alpha^{-T} \alpha^-} - y - \varepsilon.$$

$$\frac{\partial L}{\partial \alpha^+} = \cancel{K \alpha^-} - K \alpha^+ + y - \varepsilon$$

$$\frac{1}{2} K \alpha^- + \frac{1}{2} K^T \alpha^+$$

$$\frac{\partial L}{\partial \alpha^-} = \cancel{K^T \alpha^+} - K \alpha^- - y - \varepsilon.$$

$$\frac{1}{2} K \alpha^+ + \frac{1}{2} K^T \alpha^-$$

CVR