

EXP3(γ): Exponential Weights for Exploration Exploitation

Input: Time horizon T , n arms, $\eta > 0$, $\gamma \in [0, 1]$, $\lambda \in \Delta_n$.

Initialize: Player sets $p_1 = (1/n, \dots, 1/n) \in \Delta_n$. Adversary chooses $\{\ell_t\}_{t=1}^T \subset [-1, 1]^n$.

for: $t = 1, \dots, T$

Player draws $I_t \sim q_t := (1 - \gamma)p_t + \gamma\lambda$ and suffers (and observes) loss $\ell(I_t, \ell_t) = \ell_{t, I_t}$

Player computes $\hat{\ell}_{t,i} = \frac{\mathbf{1}_{\{I_t=i\}}}{q_{t,i}} \ell_{t,i}$

Update iterates:

$$w_{t+1,i} = w_{t,i} \exp(-\eta \hat{\ell}_{t,i}) \quad p_{t+1,i} = w_{t+1,i} / \sum_{j=1}^n w_{t+1,j}$$

Theorem 17. Fix any sequence $\ell_t \in [-1, 1]$ for all t and assume $\eta, \gamma \geq 0$ satisfy $-\eta \hat{\ell}_{t,i} \leq 1$ for all i, t . Then we have

$\gamma = 3d$

$$\max_{i=1, \dots, n} \mathbb{E} \left[\sum_{t=1}^T \ell_{t,i} - \ell_{t,i} \right] \leq 2\gamma T + \frac{\log(n)}{\eta} + (1 - \gamma)\eta \mathbb{E} \left[\sum_{t=1}^T \sum_{j=1}^n p_{t,j} \hat{\ell}_{t,j}^2 \right]$$

$$= 2\gamma d T + \frac{\log(n)}{\gamma} + 2d T \frac{dT}{1-\gamma} = 3\gamma d T + \frac{\log(n)}{\gamma} \leq \sqrt{12d T \log(n)}$$

Linear Bandits - Adversarial Setting

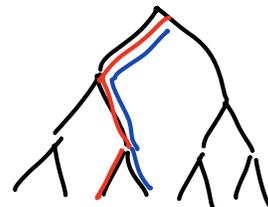
Input $\mathcal{X} = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$

Adversary chooses $\{\theta_t\}_{t=1}^T \subset \mathbb{R}^d$: $\max_{x \in \mathcal{X}} |\langle x, \theta_t \rangle| \leq 1 \quad \forall t$

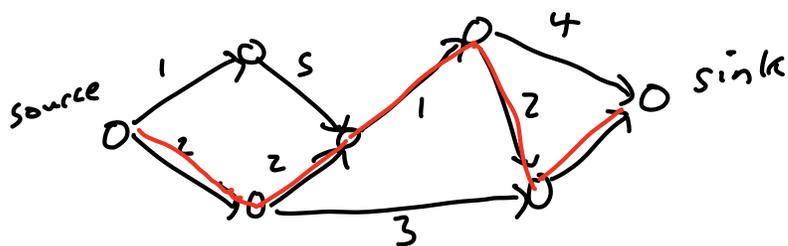
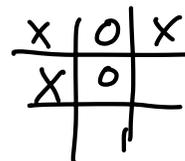
for $t = 1, 2, \dots, T$

Player chooses $x_t \in \mathcal{X}$

Adversary reveals $\langle \theta_t, x_t \rangle$ (not θ_t !)



Regret = $\max_{x \in \mathcal{X}} \sum_{t=1}^T \langle \theta_t, x_t - x \rangle$



$x \in \{0, 1\}^E$

For $x_i \in \mathcal{X}$ define $l_{\epsilon, i} = \langle x_i, \theta_\epsilon \rangle$

$$\hat{l}_{\epsilon, i} = x_i^\top A(q_\epsilon)^{-1} x_{I_\epsilon} l_{\epsilon, I_\epsilon}, \quad A(\lambda) = \sum_{i=1}^n \lambda_i x_i x_i^\top$$

$$\begin{aligned} \mathbb{E}[A(q_\epsilon)^{-1} x_{I_\epsilon} l_{\epsilon, I_\epsilon}] &= \sum_{j=1}^n q_{\epsilon, j} A(q_\epsilon)^{-1} x_j l_{\epsilon, j} \\ &= \sum_{j=1}^n q_{\epsilon, j} A(q_\epsilon)^{-1} x_j x_j^\top \theta_\epsilon \\ &= \Gamma \cdot \theta_\epsilon = \theta_\epsilon \end{aligned}$$

$$\Rightarrow \mathbb{E}[\hat{l}_{\epsilon, i}] = l_{\epsilon, i} \quad \forall i$$

Aside $x_t \sim \lambda$, $y_t = \langle x_t, \theta_0 \rangle + \underline{z}_t$ $t=1, \dots, T$

$$\hat{\theta}_{LS} = \left(\frac{1}{T} \sum_{t=1}^T x_t x_t^\top \right)^{-1} \frac{1}{T} \sum_{t=1}^T x_t y_t$$

$$\mathbb{E}[\hat{\theta}_{LS}] = \theta_0, \quad \text{Cov}(\hat{\theta}_{LS}) = \left(\sum_{t=1}^T x_t x_t^\top \right)^{-1}$$

$$\hat{\theta}_{IPS} = \left(\sum_{x \in \mathcal{X}} \lambda_x x x^\top \right)^{-1} \frac{1}{T} \sum_{t=1}^T x_t y_t$$

$$\mathbb{E}[\hat{\theta}_{IPS}] = \theta_0, \quad \text{Cov}(\hat{\theta}_{IPS}) = \left(T \sum_{x \in \mathcal{X}} \lambda_x x x^\top \right)^{-1}$$

$$x = \{[1], [0]\}$$

$$\lambda_1, \lambda_2 = 1 - \lambda,$$

$$T_1 = \sum_{t=1}^T \mathbb{1}\{I_t = 1\}$$

$$\hat{\theta}_{LS} = \begin{bmatrix} \frac{1}{T_1} \sum_{t=1}^T \mathbb{1}\{I_t = 1\} y_t \\ \frac{1}{T - T_1} \sum_{t=1}^T \mathbb{1}\{I_t = 2\} y_t \end{bmatrix}$$

$$\hat{\theta}_{IPS} = \begin{bmatrix} \frac{1}{\lambda_1 T} \sum_{t=1}^T \mathbb{1}\{I_t = 1\} y_t \\ \frac{1}{(1-\lambda_1) T} \sum_{t=1}^T \mathbb{1}\{I_t = 1\} y_t \end{bmatrix}$$

$$\hat{l}_{t,i} = x_i^T A(q_t)^{-1} x_{I_t} l_{t,I_t}$$

$$- \sum \hat{l}_{t,i} \leq \sum |x_i^T A(q_t)^{-1} x_{I_t}| \quad (|l_{t,i}| \leq 1)$$

$$= \sum \left| \langle A(q_t)^{-1/2} x_i, A(q_t)^{-1/2} x_{I_t} \rangle \right|$$

$$\leq \sum \|x_i\|_{A(q_t)^{-1}} \cdot \|x_{I_t}\|_{A(q_t)^{-1}}$$

$$\leq \max_{j=1, \dots, n} \sum \|x_j\|_{A(q_t)^{-1}}^2$$

$$\leq \max_{j=1, \dots, n} \frac{\sum \|x_j\|_{A(\lambda)^{-1}}^2}{\delta}$$

$$= \frac{\sum}{\delta} \cdot d$$

if λ δ -optimal.

Suggests picking $\gamma = \alpha d$ s.t. $-\exists \hat{l}_{t,i} \leq 1 \quad \forall i, t.$

$$\begin{aligned} \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^n P_{t,i} \hat{l}_{t,i}^2 \right] &= \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^n P_{t,i} x_i^T A(\rho_t)^{-1} x_{\mathcal{I}_t} x_{\mathcal{I}_t}^T A(\rho_t)^{-1} x_i l_{t,\mathcal{I}_t}^2 \right] \\ &\leq \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^n P_{t,i} x_i^T A(\rho_t)^{-1} x_{\mathcal{I}_t} x_{\mathcal{I}_t}^T A(\rho_t)^{-1} x_i \right] \\ &= \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^n P_{t,i} x_i^T A(\rho_t)^{-1} x_i \right] \\ &\leq \frac{1}{1-\gamma} \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^n P_{t,i} x_i^T A(\rho_t)^{-1} x_i \right] \\ &= \frac{1}{1-\gamma} \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^n P_{t,i} \text{Trace} \left(x_i^T A(\rho_t)^{-1} x_i \right) \right] \\ &= \frac{1}{1-\gamma} \mathbb{E} \left[\sum_{t=1}^T \sum_i P_{t,i} \text{Trace} \left(A(\rho_t)^{-1} x_i x_i^T \right) \right] \\ &= \frac{dT}{1-\gamma} \end{aligned}$$

Contextual Bandits

for $t=1, 2, \dots, T$

Nature reveals context $c_t \stackrel{iid}{\sim} \mathcal{V}$

Player chooses $I_t \in [n]$

and gets reward $y_t \in [-1, 1]$, s.t. $E[y_t | c_t] = r(c_t, I_t)$

Suppose $\text{support}(\mathcal{V}) = \mathcal{C}$ and $|\mathcal{C}| < \infty$.

Treat each context as indep. bandit.

$$\sum_{t=1}^T \max_i (r(c_t, i) - r(c_t, I_t))$$

$$= \sum_{c \in \mathcal{C}} \max_i \sum_{t=1}^T \mathbb{1}\{c_t = c\} (r(c_t, i) - r(c_t, I_t))$$

$$\leq \sum_{c \in \mathcal{C}} \sqrt{n T_c}$$

← If playing MAB w/ \sqrt{nT} regret

$$T_c = \sum_t \mathbb{1}\{c_t = c\}$$

$$\leq \sqrt{\left(\sum_{c \in \mathcal{C}} \sqrt{T_c}\right)^2 \cdot \left(\sum_{c \in \mathcal{C}} \sqrt{n}\right)^2}$$

$$= \sqrt{n |\mathcal{C}| T}$$

Alternatively, what if we ignore context?

$$\max_i \sum_{t=1}^T r(c_t, i) - r(c_t, I_t) \leq \sqrt{nT}$$

Total Reward:

Bandit per context $\sum_{t=1}^T r(c_t, I_t) \geq \sum_{t=1}^T \max_i r(c_t, i) - \sqrt{n|C|T}$

$$|\Pi| = n^{|C|}$$

Ignore context $\sum_{t=1}^T r(c_t, I_t) \geq \max_i \sum_{t=1}^T r(c_t, i) - \sqrt{nT}$

$$|\Pi| = n$$

Policy regret $\sum_{t=1}^T r(c_t, a_t) \geq \max_{\pi \in \Pi} \sum_{t=1}^T \mathbb{E}[r(c_t, i)] - \sqrt{nT \log |\Pi|}$

Consider a finite set of policies Π , $\pi \in \Pi$ $\pi: C \rightarrow [n]$

Policy regret $\max_{\pi \in \Pi} \sum_{t=1}^T r(c_t, \pi) - r(c_t, \hat{\pi}_t)$

———— $\pi(c_t) \in \Delta_d$, $a_t \sim \pi(c_t)$ then $P(j = a_t | c_t) = \pi(j | c_t)$

Suppose contexts are arbitrary (adversarial!!)

Adversary chooses $\{c_t\}_{t=1}^T$ in advance

for $t=1, 2, \dots, T$

Adversary reveals c_t

Player: chooses $I_t \in [n]$, draws $a_t \in \pi(c_t)$

gets loss $l(c_t, a_t)$

Regret $\max_{\pi \in \Pi} \sum_{t=1}^T \mathbb{E}_{a_t \sim \pi(c_t)} [l(c_t, a_t)] - \mathbb{E}_{a_t \sim \pi(c_t)} [l(c_t, a)]$

EXP3(γ): Exponential Weights for Exploration Exploitation

Input: Time horizon T , n arms, $\eta > 0$, $\gamma \in [0, 1]$

Initialize: Player sets $p_1 = (1/n, \dots, 1/n) \in \Delta_n$. Adversary chooses $\{\ell_t\}_{t=1}^T \subset [-1, 1]^n$.

for: $t = 1, \dots, T$ *Adversary reveals c_t*

Player defines $\lambda_t \in \Delta_n$ and plays $I_t \sim q_t := (1 - \gamma)p_t + \gamma\lambda_t$

Player suffers (and observes) loss ℓ_{t, I_t} (but does *not* observe $\ell_{t, i}$ for $i \neq I_t$)

Player computes loss estimator $\hat{\ell}_{t, i}$ for all $i \in [n]$

Update iterates:

$$w_{t+1, i} = w_{t, i} \exp(-\eta \hat{\ell}_{t, i}) \quad p_{t+1, i} = w_{t+1, i} / \sum_{j=1}^n w_{t+1, j}$$

Theorem 17. Fix any sequence $\ell_t \in [-1, 1]$ for all t . For any $\hat{\ell}_{t, i}$ and $\eta, \gamma \geq 0$ that satisfy $\mathbb{E}[\hat{\ell}_{t, i} | \mathcal{F}_{t-1}] = \ell_{t, i}$ and $-\eta \hat{\ell}_{t, i} \leq 1$ for all i, t we have

$$\gamma = \frac{3}{2}$$

d -actions
 n -policies

$$\max_{i=1, \dots, n} \mathbb{E} \left[\sum_{t=1}^T \ell_{t, I_t} - \ell_{t, i} \right] \leq 2\gamma T + \frac{\log(n)}{\eta} + (1 - \gamma)\eta \mathbb{E} \left[\sum_{t=1}^T \sum_{j=1}^n p_{t, j} \hat{\ell}_{t, j}^2 \right]$$

$$\ell_{t, i} := \mathbb{E}_{a_t \sim \pi_i} [\ell(c_t, a_t)]$$

$$\leq 3 \frac{3}{2} T + \frac{\log(n)}{\frac{3}{2}}$$

$$\text{I observe } \ell(c_t, a_t)$$

$$\leq \sqrt{12 d T \log(n)}$$

$$\hat{\ell}_{t, i} = \frac{\pi_i(a_t | c_t)}{\sum_{k=1}^n q_{t, k} \pi_k(a_t | c_t)} \ell(c_t, a_t)$$

$$a_t \sim \pi_{I_t}(c_t)$$

$$I_t \sim q_t$$

$$\mathbb{E}[\hat{\ell}_{t, i}] = \mathbb{E} \left[\frac{\pi_i(a_t | c_t) \ell(c_t, a_t)}{\sum_k q_{t, k} \pi_k(a_t | c_t)} \right]$$

$$= \sum_{j=1}^d \mathbb{E} \left[\underbrace{P(a_t = j)} \frac{\pi_i(j | c_t) \ell(c_t, j)}{\sum_k q_{t,k} \pi_k(j | c_t)} \right]$$

$$= \sum_{s=1}^n P(a_t = j | I_t = s) P(I_t = s)$$

$$= \sum_{s=1}^n \pi_s(j | c_t) q_{t,s}$$

$$= \sum_{s=1}^d \mathbb{E} \left[\pi_i(j | c_t) \ell(c_t, j) \right]$$

$$= \mathbb{E}_{a \sim \pi(c_t)} \left[\ell(c_t, a) \right] = \ell_{t,i}$$

$$- \hat{\ell}_{t,i} \leq \sum_{k=1}^n \frac{\pi_i(a_t | c_t)}{\sum_{k=1}^n q_{t,k} \pi_k(a_t | c_t)}$$

$$\leq \sum_{j=1}^d \frac{\pi_i(j | c_t)}{\sum_{k=1}^d q_{t,k} \pi_k(a_t | c_t)}$$

$$\leq \frac{\sum_{i=1}^d \alpha_i(j(\epsilon))}{\sum_{n=1}^n \lambda_{\epsilon, n} \tau_n(u_\epsilon(\epsilon))}$$

$$\leq \frac{\sum d}{\gamma} \leq 1 \quad \text{if } \gamma = \sum d.$$

Proposition For $i=1, \dots, n$ let $p_i \in \Delta_d$. Then

$$\min_{\lambda \in \Delta_n} \max_{i=1, \dots, n} \sum_{j=1}^d \frac{p_{i,j}}{\sum_k \lambda_k p_{k,j}} = d.$$

Proof $f(\lambda) := \sum_{j=1}^d \log \left(\sum_{k=1}^n \lambda_k p_{k,j} \right), \quad \lambda^* = \underset{\lambda}{\text{argmax}} f(\lambda)$

for any λ $f(\lambda) \leq f(\lambda^*) + \nabla f(\lambda^*)^T (\lambda - \lambda^*)$

$$\begin{aligned} \Rightarrow 0 &\geq \nabla f(\lambda^*)^T (\lambda - \lambda^*) \\ &= \sum_{i=1}^n \sum_{j=1}^d \frac{p_{i,j}}{\sum_{k=1}^n p_{k,j} \lambda_k^*} (\lambda_i - \lambda_i^*) \\ &= \sum_{i=1}^n \sum_{j=1}^d \frac{p_{i,j} \lambda_i}{\sum_{k=1}^n p_{k,j} \lambda_k^*} - d \end{aligned}$$

Take $\lambda = e_i$

$$\Rightarrow \sum_{j=1}^d \frac{P_{i,j}}{\sum_{k=1}^n P_{k,i} \lambda_k} \leq d.$$

$$\min_{\lambda \in \Delta_n} \max_{i=1, \dots, n} \sum_{j=1}^d \frac{P_{i,j}}{\sum_k \lambda_k P_{k,i}} \leq \max_{i=1, \dots, n} \sum_{j=1}^d \frac{P_{i,j}}{\sum_k \lambda_k P_{k,i}} \leq d$$

$$\min_{\lambda \in \Delta_n} \max_{i=1, \dots, n} \sum_{j=1}^d \frac{P_{i,j}}{\sum_k \lambda_k P_{k,i}} \geq \min_{\lambda} \sum_{i=1}^n \lambda_i \sum_{j=1}^d \frac{P_{i,j}}{\sum_k \lambda_k P_{k,i}} = d$$

□

$$\mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^n P_{e_t, i} \hat{\ell}_{e_t, i}^2 \right] \leq \frac{dT}{1-\gamma}$$

Input policy set Π , $\pi \in \Pi$, $\pi: \mathcal{C} \rightarrow \mathcal{X}$

for $t=1, 2, \dots, T$

Nature reveals $(c_t^{\text{oil}}) \vee$

Player chooses $x_t \in \mathcal{X}$ observes $r(c_t, x_t) + \gamma_t$

Collection of policies Π , $\pi: \mathcal{C} \rightarrow \mathcal{X}$

Off-policy Estimation

Suppose we play a logging policy $\mu(\cdot | c_t) \in \Delta_{\mathcal{X}}$

$x_t \sim \mu(\cdot | c_t) \quad \forall t$ observe $y_t = r(c_t, x_t) + \gamma_t$

We get dataset $\left\{ (c_t, x_t, y_t, p_t) \right\}_{t=1}^T$

$$p_t = \mu(x_t | c_t)$$

$$V(\pi) = \mathbb{E}_{c_t \sim \nu} [r(c_t, \pi(c_t))]$$

$$\hat{V}(\pi) = \frac{1}{T} \sum_{t=1}^T \frac{\mathbb{1}\{x_t = \pi(c_t)\}}{p_t} y_t$$

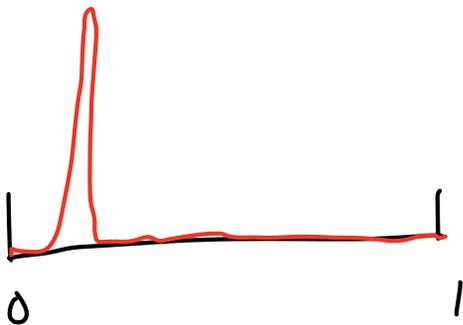
$$\mathbb{E}[\hat{V}(\pi)] = V(\pi)$$

$$\begin{aligned} \mathbb{E}[(\hat{V}(\pi) - V(\pi))^2] &= \mathbb{E}\left[\mathbb{E}[(\hat{V}(\pi) - V(\pi))^2 | \{\xi_t\}_{t=1}^T}\right] \\ &\leq \mathbb{E}\left[\frac{1}{T^2} \sum_{t=1}^T \frac{1}{P_t}\right] \\ &\leq \frac{1}{T} \mathbb{E}_{c \sim \nu} \left[\frac{1}{\mu(\pi(c) | c)}\right]. \end{aligned}$$

Bernstein's Inequality | Let X_1, \dots, X_m be iid R.V. w/ $\frac{1}{m} \sum_{i=1}^m \mathbb{E}[(X_i - \mathbb{E}[X_i])^2] \leq \sigma^2$ and $|X_i| \leq B$ then

$$\left| \frac{1}{m} \sum_{i=1}^m X_i - \mathbb{E}[X_i] \right| \leq \sqrt{\frac{2\sigma^2 \log(2/\delta)}{m}} + \frac{2B \log(2/\delta)}{3m}.$$

w.p. $\geq 1 - \delta$.



Corollary $|\hat{V}(\pi) - V(\pi)| \leq \underbrace{\sqrt{\frac{2 \log(2/\pi/\delta) \cdot \mathbb{E}\left[\frac{1}{\mu(\pi(c_t)|c_t)}\right]}{T}}}_{\bar{V}} + \frac{2\bar{V} \log(2/\pi/\delta)}{3T}$

w.p. $\geq 1-\delta$ where $\bar{V} = \max_{c, x} \frac{1}{\mu(x|c)} \cdot C_T(\pi)$

Uniform explanation

$$\mu(x|c) = \frac{1}{|\mathcal{X}|} \Rightarrow \mathbb{E}\left[\frac{1}{\mu(\cdot|\cdot)}\right] = |\mathcal{X}|$$

$$\bar{V} = |\mathcal{X}|$$

$$\begin{aligned} \Rightarrow |\hat{V}(\pi) - V(\pi)| &\leq \sqrt{\frac{2 \log(2/\pi/\delta) |\mathcal{X}|}{T}} + \frac{2|\mathcal{X}| \log(2/\pi/\delta)}{3T} \\ &\leq \sqrt{\frac{4 |\mathcal{X}| \log(2/\pi/\delta)}{T}} \end{aligned}$$

$$\hat{\pi}_{MLE} = \underset{\pi \in \Pi}{\operatorname{argmax}} \hat{V}(\pi)$$

$$= \underset{\pi \in \Pi}{\operatorname{argmax}} \frac{1}{T} \sum_{t=1}^T \frac{\mathbb{1}\{\pi(c_t) = x_t\}}{P_t} y_t$$

$$= \underset{\pi}{\operatorname{argmax}} \frac{1}{T} \sum_{t=1}^T (1 - \mathbb{1}\{\pi(c_t) \neq x_t\}) \frac{y_t}{P_t}$$

$$= \underset{\pi \in \Pi}{\operatorname{argmin}} \frac{1}{T} \sum_{t=1}^T \mathbb{1}\{\pi(c_t) \neq x_t\} \frac{y_t}{P_t}$$

Consider generic machine learning problem $\{(x_i, y_i)\}_{i=1}^n$,
 $x_i \in \mathbb{R}^d$, $y_i \in [k]$. Find a classifier that predicts

y_i from x_i . Equiv. find a classifier that minimizes
 misclassification $\frac{1}{n} \sum_{i=1}^n \mathbb{1}\{h(x_i) \neq y_i\}$.

$$\approx \frac{1}{n} \sum_{i=1}^n -\log \left(\frac{\exp([f_\theta(x_i)]_{y_i})}{\sum_k \exp([f_\theta(x_i)]_k)} \right)$$

$$f_\theta: \mathbb{R}^d \rightarrow \mathbb{R}^k$$

Predicting $P(y_i = j | x_i) = \frac{\exp([f_\theta(x_i)]_j)}{\sum_k \exp([f_\theta(x_i)]_k)}$

$$V(\hat{\tau}_{MLE}) = \hat{V}(\hat{\tau}_{MLE}) + V(\hat{\tau}_{MLE}) - \hat{V}(\hat{\tau}_{MLE})$$

$$\geq \hat{V}(\hat{\tau}_{MLE}) - C_T(\hat{\tau}_{MLE})$$

$$\geq \hat{V}(\tau_\alpha) - C_T(\hat{\tau}_{MLE})$$

$$\geq V(\bar{\pi}_a) - C_T(\bar{\pi}_a) - C_T(\hat{\pi}_{MLE})$$

$$\geq V(\bar{\pi}_a) - \max_{\pi \in \Pi} C(\pi)$$

Principle of Pessimism

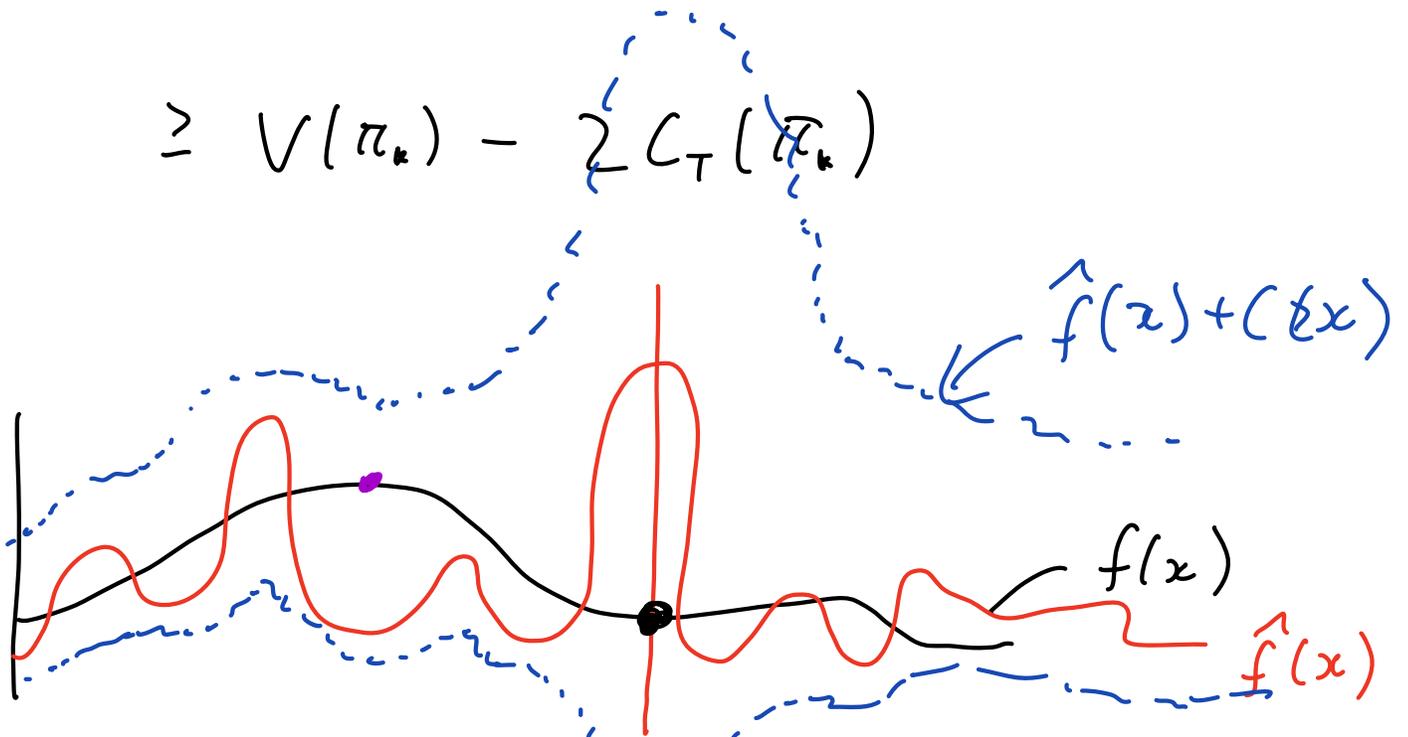
$$\hat{\pi}_{pes} = \underset{\pi \in \Pi}{\operatorname{argmax}} \hat{V}(\pi) - C_T(\pi)$$

$$V(\hat{\pi}_{pes}) = \hat{V}(\hat{\pi}_{pes}) + V(\hat{\pi}_{pes}) - \hat{V}(\hat{\pi}_{pes})$$

$$\geq \hat{V}(\hat{\pi}_{pes}) - C_T(\hat{\pi}_{pes})$$

$$\geq \hat{V}(\pi_*) - C_T(\pi_*)$$

$$\geq V(\pi_*) - C_T(\pi_*)$$



x

$\uparrow \hat{f}(x) - c(x)$

$$\hat{x} = \underset{x}{\operatorname{argmax}} \hat{f}(x)$$

$$\mathbb{E}[f(\hat{x})] \leq \max_x f(x)$$

$$\mathbb{E}[\hat{f}(\hat{x})] > \max_x f(x)$$

