

- Supervised Learning
- Unsupervised Learning
- Semi-supervised learning

\Rightarrow In supervised learning, we have loss functions such as l_{0-1} and corresponding risks

$$R_{0-1} = E[l_{0-1}(f(x), y)]$$

* Bayes Classifier: \rightarrow Pointwise minimizer
(Ideal classifier) (For 2 class classification ± 1)

$$\eta(x) = P(y=1|x)$$

$$f^*(x) = \text{sign}(\eta(x) - 1/2) \rightarrow \text{Bayes classifier}$$

Show,

$$\Rightarrow R_{0-1}(g) \geq R_{0-1}(f^*) \quad \forall g \neq f^*$$

$$R_{0-1}(g) = E[\mathbb{I}_{\{y \cdot g(x) < 0\}}]$$

$= P(y \cdot g(x) < 0) \rightarrow$ Probability of misclassification

$$\bullet \text{ Show, } R_{0-1}(g) - R_{0-1}(f^*) \geq 0$$

$$E[\mathbb{I}_{\{y \neq g(x)\}} - \mathbb{I}_{\{y \neq f^*(x)\}}] \geq 0$$

$$\Rightarrow E_x [E_{y|x} [\mathbb{I}_{\{y \neq g(x)\}}|x] - E_{y|x} [\mathbb{I}_{\{y \neq f^*(x)\}}|x]] \geq 0$$

$$\bullet E_{y|x} [\mathbb{I}_{\{y \neq g(x)\}}|x] = P[y \neq g(x)|x] \\ = 1 - P[y = g(x)|x]$$

$$= 1 - P[y=1, g(x)=1|x]$$

$$= 1 - \mathbb{E} [\mathbb{I}_{\{y=1\}} \mathbb{I}_{\{g(x)=1\}} | x]$$

$$- \mathbb{E} [\mathbb{I}_{\{y=-1\}} \mathbb{I}_{\{g(x)=-1\}} | x]$$

$$= 1 - \mathbb{I}_{\{g(x)=1\}} \mathbb{E} [\mathbb{I}_{\{y=1\}} | x]$$

$$- \mathbb{I}_{\{g(x)=-1\}} \mathbb{E} [\mathbb{I}_{\{y=-1\}} | x]$$

Now,

$$\textcircled{1} \leftarrow \mathbb{E}_{y|x} [\mathbb{I}_{\{y \neq g(x)\}} | x] = 1 - \mathbb{I}_{\{g(x)=1\}} \cdot h(x) - \mathbb{I}_{\{g(x)=-1\}} (1-h(x))$$

$$\textcircled{2} \leftarrow \mathbb{E}_{y|x} [\mathbb{I}_{\{y \neq f^*(x)\}} | x] = 1 - \mathbb{I}_{\{f^*(x)=1\}} \cdot h(x)$$

$$\textcircled{1} - \textcircled{2} = h(x) [\mathbb{I}_{\{f^*(x)=1\}} - \mathbb{I}_{\{g(x)=1\}}] + (1-h(x)) [\mathbb{I}_{\{f^*(x)=-1\}} - \mathbb{I}_{\{g(x)=-1\}}]$$

$$\textcircled{1} - \textcircled{2} = h(x) [\mathbb{I}_{\{f^*(x)=1\}} - \mathbb{I}_{\{g(x)=1\}}] + (h(x)-1) [\mathbb{I}_{\{f^*(x)=1\}} - \mathbb{I}_{\{g(x)=1\}}]$$

$$\textcircled{1} - \textcircled{2} = (2h(x)-1) [\mathbb{I}_{\{f^*(x)=1\}} - \mathbb{I}_{\{g(x)=1\}}]$$

Case ① : $\eta(x) > 1/2$, $2\eta(x) - 1 > 0$

Probability that $y=1$ given x .

$$\mathbb{I}_{\{g^*(x)=1\}} = 1 \quad (\because \eta(x) > 1/2)$$

$$\therefore ① - ② \geq 0$$

$$\therefore \mathbb{E}_x [(2\eta(x) - 1) [\mathbb{I}_{\{g^*(x)=1\}} - \mathbb{I}_{\{g(x)=1\}}]] \geq 0$$

Case ② : $\eta(x) < 1/2 \Rightarrow 2\eta(x) - 1 < 0$

$$\mathbb{I}_{\{g^*(x)=1\}} = 0$$

$$\therefore \mathbb{I}_{\{g^*(x)=1\}} - \mathbb{I}_{\{g(x)=1\}} \leq 0$$

$$\therefore \mathbb{E}_x [(2\eta(x) - 1) [\mathbb{I}_{\{g^*(x)=1\}} - \mathbb{I}_{\{g(x)=1\}}]] \geq 0$$

$$\therefore R_{0-1}(g) - R_{0-1}(g^*) \geq 0$$

- All standard loss functions like log, exponential etc are classification calibrated (Achieve bayes like classifiers)

- Empirical vs Expected risk \rightarrow Large number of samples required to reduce gap between the two.

(2)

Note! Underlying distribution is fixed but unknown in any ML related data source.

(We use Empirical risk (average) since we don't know the true expected risk)

Online Learning: (One sample at a time)

- Stochastic Gradient Descent (to update the hypothesis)

$$w_{t+1} = w_t - \eta (\nabla_{w_t} L(f(x_t), y_t))$$

(For each datapoint)

*1) Perceptron: (Single sample algorithm)

$$\bullet g = \text{sign}(f(x)) = \text{sign}\left(\sum_{i=1}^d w_i x_i + b\right)$$

$$x = [x_1 \ x_2 \ \dots \ x_d]^T$$

$$\bullet \text{Loss: } L_{\text{perceptron}}(f(x), y) = \max(0, -y \cdot f(x))$$

(Not a 0-1 loss function)

$$\bullet w_{t+1} = \begin{cases} w_t + \eta y_t x_t & , y_t (w^T x_t + b) \leq 0 \\ w_t & , y_t (w^T x_t + b) > 0 \end{cases}$$

$$\bullet b_{t+1} = \begin{cases} b_t + \eta y_t & , y_t (w^T x_t + b) \leq 0 \\ b_t & , y_t (w^T x_t + b) > 0 \end{cases}$$

Show,

$$\Rightarrow y_t (w^T x_t + b_t) < y_t (w_{t+1}^T x_t + b_{t+1})$$

(After an update, the point is classified better)

$$= y_t (w_t^T x_t +$$

$$+ \eta y_t \|x_t\|^2)$$

$$= y_t (w_t^T x_t + b_t) + \eta (||x_t||^2 + 1) \cdot y_t + b_t + \eta y_t \rightarrow \text{the value.}$$

Theorem 1: Finite convergence in linearly separable case.

* * * Let $S = \{(x_1, y_1), \dots, (x_T, y_T)\}$ be the examples in T iterations

(Assume bias=0, can be done)

$$\text{by } \tilde{\omega} = \begin{bmatrix} \omega \\ b \end{bmatrix}$$

$$\tilde{x} = \begin{bmatrix} x \\ 1 \end{bmatrix}$$

Ideal classifier (Assume)

- Let $\exists u \in \mathbb{R}^d$ s.t $\|u\|_2 = 1$ margin (Ideal)
- and $y_t \cdot u \cdot x_t \geq \gamma \forall t = 1, \dots, T \text{ & } \gamma > 0$

$$\text{Let } R_2 = \max_{\tilde{x}} \|\tilde{x}\|$$

Radius from which examples are taken.

\Rightarrow Perceptron converges in $\left(\frac{R_2^2}{\gamma^2}\right) \text{ misclassified examples}$.

Proof: Let x_t gets misclassified, ($u = \text{Ideal classifier}$)

$$y_t (\cancel{u \cdot x_t}) < 0$$

$$w_{t+1} \cdot u = (w_t + y_t x_t) \cdot u$$

$$= w_t \cdot u + \underbrace{y_t x_t \cdot u}_{\geq \gamma}$$

$$\geq w_t \cdot u + \gamma$$

$$w_{t+1} \cdot u - w_t \cdot u \geq \gamma$$

- Suppose in T trials, k datapoints are misclassified

$$w_{T+1} \cdot u \geq k\gamma \quad (w_1 = 0)$$

$a \cdot b \leq \|a\| \|b\| \quad \} \text{ Cauchy-Schwarz inequality}$

$$\Rightarrow w_{t+1} \cdot u \leq \|u\| \cdot \|w_{t+1}\| = \|w_{t+1}\|$$

$$\Rightarrow \|w_{t+1}\|^2 = \|w_t + y_t x_t\|^2$$

$$= \|w_t\|^2 + \|x_t\|^2 + 2y_t w_t \cdot x_t$$

$$\leq \|w_t\|^2 + R_2^2$$

$$\therefore \|w_{t+1}\|^2 \leq \|w_t\|^2 + k \cdot R_2^2$$

$$\text{So } \|w_{t+1}\|^2 \leq k \cdot R_2^2 \quad (\text{k mistakes in T iterations})$$

$$\text{Now, } \|w_{t+1}\| \geq w_{t+1} \cdot u \geq k\gamma$$

$$\Rightarrow k\gamma \leq \|w_{t+1}\|$$

$$k^2 \gamma^2 \leq \|w_{t+1}\|^2 \leq kR_2^2$$

$$\Rightarrow k \leq \frac{R_2^2}{\gamma^2}$$

At most k misclassified for convergence.

Note: Perceptron extension proofs (Multiclass etc)

* Note: $k = \frac{R_2^2}{\gamma^2}$ for examples like $x_i = (0, 0, \dots, 1, 0)$

and u_i here would be $\frac{y_i}{\|u_i\|}$

$$R_2^2 = 1, \gamma = \frac{1}{\|u_i\|}$$

(Here $k = n$ iterations)

Final value of $\|w\|$ (Final value)

Theorem 2: $S = \{(x_1, y_1), \dots, (x_T, y_T)\}$ (Finite convergence in non linearly separable case)

** $(x_i, y_i) \in \mathbb{R}^d \times \{-1, 1\}$ $\forall i$

$$\bullet R_2 = \max_{\pm} \|x_i\|_2, \text{ let } \gamma > 0$$

Then for any $u \in \mathbb{R}^d (\|u\|_2 = 1)$

the number of mistakes

$$R \leq \frac{(R_2 + \delta)^2}{\gamma^2}$$

$$\Rightarrow D^2 = \sum_{t=1}^T d_t^2 \text{ and } d_t = \max_{y - u \cdot x_t} \delta \quad (\text{Allowed error})$$

- Classifier won't be able to classify perfectly in current space, so we project data.

$x_t \rightarrow x_t'$ where,

$$x_t' = [x_{t,1}, x_{t,2}, \dots, x_{t,d}, 0, \dots, 0]$$

\downarrow
d+1 dimension
 $(d+1)^{\text{th}}$ component

$$\bullet u' = \left[\frac{u_1}{z}, \frac{u_2}{z}, \dots, \frac{u_d}{z}, \frac{y_1 d_1}{z}, \dots, \frac{y_T d_T}{z} \right]$$

Set, $z = \sqrt{1 + \frac{D^2}{\Delta^2}}$ when we equate $\|u'\|_2 = 1$

$$\bullet y_t u' x_t' = y_t \left[\frac{u \cdot x_t}{z} + \frac{\Delta \cdot y_t d_t}{z} \right]$$

$$= \frac{y_t \cdot u \cdot x_t}{z} + \frac{d_t}{z}$$

$$\geq \frac{y_t \cdot u \cdot x_t}{z} + \frac{\gamma - y_t \cdot u \cdot x_t}{z}$$

$$\therefore 0 = y_t u' x_t \rightarrow y_t/z \rightarrow 1$$

$$\Rightarrow \|x + \zeta\|^2 = \|x + \zeta\|^2 + \Delta^2 \leq R_2^2 + \Delta^2$$

(Using ① & ②,
similar to
Theorem 1,
we can get
the two
inequalities
comparing
which we get
k)

$$k \leq \frac{(R_2^2 + \Delta^2) \cdot z^2}{\gamma^2}, \text{ Substituting } z$$

$$k \leq \frac{(R_2^2 + \Delta^2)(1 + \Delta^2/R_2^2)}{\gamma^2}$$

! () $(\Delta = \sqrt{R_2 D}) \rightarrow$ (To minimize k, differentiate w.r.t. Δ)

$$(\because w_{++}^\top u - w_+^\top u \geq \gamma/z)$$

$$(\because \|w_{++}\|^2 = \|w_+\|^2 + R_2^2 + \Delta^2)$$

③

*SVM Recap:

$$w = \sum_{i=1}^N \alpha_i y_i x_i$$

$$f(x) = w^\top x + b$$

$$= \sum_{i=1}^N \alpha_i y_i (x_i^\top x) + b$$

Primal: (with slack)

$$\min_{w, b, \xi} \frac{1}{2} \|w\|_2^2 + C \cdot \sum_{i=1}^N \xi_i$$

$$\text{s.t. } y_i(w^\top x_i + b) \geq 1 - \xi_i \quad \forall i$$

$$\xi_i \geq 0 \quad \forall i$$

↓

Dual: (with slack)

$$\max_{\alpha_1, \dots, \alpha_N} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j (x_i^\top x_j)$$

$$\begin{cases} \bar{x}^\top \mathbf{1} - \frac{1}{2} (\bar{x} \circ \bar{y})^\top \\ G(\bar{x} \circ \bar{y}) \end{cases} \quad \text{s.t. } 0 \leq \alpha_i \leq C \quad \forall i$$

$$\sum_{i=1}^N \alpha_i y_i = 0$$

- Replace $x_i^T x_j$ with $\phi(x_i, x_j)$, some kernel function to project non-linearly separable data onto higher dimensions.

$$\overbrace{\hspace{1cm}}^x$$

\Rightarrow Perceptron:

(details later) $w^T = 0$ $f(x) = w^T x + b$

$w_i = \sum_{j \in S} y_j x_{ij}$

where $S = \{j \in \{1, \dots, d\} \mid y_j w_j \leq 0\}$

- This is similar to the SVM and we can use the kernel trick here as well.

Kernel function:

$$K(x_1, x_2) = \phi(x_1)^T \cdot \phi(x_2)$$

$$\overbrace{\hspace{1cm}}^x$$

*Kernels:

- Linear kernel $\Rightarrow K(x_1, x_2) = x_1^T x_2$

- Polynomial kernel $\Rightarrow K(x_1, x_2) = (x_1^T x_2 + 1)^d$

ex: $K(x_1, x_2) = (x_{11} x_{21} + x_{12} x_{22} + 1)^2$

$$= \phi(x_1)^T \cdot \phi(x_2)$$

where $\phi(x_i) = [x_{i1}^2, x_{i2}^2, \sqrt{2}x_{i1},$

$$\sqrt{2}x_{i1}x_{i2}, \dots, \sqrt{2}x_{i1}x_{i2}, \dots, \sqrt{2}x_{i1}x_{i2}]$$

- Gaussian Kernel $\Rightarrow K(x_1, x_2) = \exp\left(-\frac{\|x_1 - x_2\|^2}{2\sigma^2}\right)$

Projects to ∞ dimensions

(These kernels are all symmetric)

Gram Matrix: (A similarity matrix for given data)

- $G_{ij} = K(x_i, x_j)$ for a given set of vectors x_1, x_2, \dots, x_n

→ Positive Semidefinite Kernels:

- A symmetric kernel function is finite p.s.d if gram matrices formed by any finite subset of X is positive semi definite.

* Mercer's Theorem:

- If kernel function is p.s.d there exists a map $\phi: X \rightarrow H$ ^{Hilbert space} s.t $K(x_1, x_2) = \phi(x_1)^T \cdot \phi(x_2)$

XX

* Kernel based Perceptrons:

$$\bullet w = \sum_{i \in H} y_i x_i \quad \text{where } H = \{t \in \{1, \dots, T\} \mid$$

(Assume no bias, $w \leftarrow [w_b], x \leftarrow [x]$) $y_t w + x_t \leq 0$

$$\Rightarrow w = \sum_{k=1}^n \alpha_k \cdot y_k \cdot \phi(x_k)$$

Weight for training sample

Update function:

(Initially all α_k 's are zero)

- i-th sample { • If $\text{sign}(w^T \phi(x_i)) \neq y_i$,

misclassified

$x_i \rightarrow x_{i+1}$

$$w = w + y_i \phi(x_i)$$

$$\therefore w = \sum_{k \neq i} \alpha_k y_k \phi(x_k) + (\alpha_{i+1}) y_i \phi(x_i)$$

⇒ Essentially, if $\text{sign}(\sum_{k=1}^n \alpha_k y_k \phi(x_k)^T \cdot \phi(x_i)) \neq y_i$

\Rightarrow For ($t=1$ to T) \rightarrow Choose sample (x_{t+1}, y_{t+1}) randomly.
 If $y_t \left(\sum_{k=1}^{t-1} \alpha_k y_k K(x_k, x_t) \right) \leq 0$
 $\quad \quad \quad \alpha_t' \rightarrow \alpha_t + 1$

if $y_t \left(\sum_{k=1}^{t-1} \alpha_k y_k K(x_k, x_t) \right) > 0$
 else

$$\alpha_t' \rightarrow \alpha_t$$

\Rightarrow Final classifier, $f(x) = \sum_{t=1}^T \alpha_t y_t K(x_t, x)$

Convergence of Kernel Perceptron:

i) Separable case:

$$y_t K(u, x_t) \geq \gamma \quad \forall t$$

$$R_2^2 = \max_{t=1-T}^T K(x_t, x_t), \|\phi(u)\|_2^2 = 1$$

(At most $\frac{R_2^2}{\gamma^2}$ mistakes)

(Proof similar to non-kernel version)

Note: Perception vs SVM, Perception just finds any classifier which separates train data well but SVM tries to maximize the margin. (We can start including margins in perception as well)

Note: Passive ~~aggressive~~ ^{aggressive} perception algorithms
 (2006 JMLR Online Passive Aggressive)

Margin $\geq \gamma$ to be

④ Learning from experts:

$$R_T = \sum_{t=1}^T L(\hat{y}_t, y_t) - \min_{i=1 \dots h} \sum_{t=1}^T L(\hat{y}_t^i, y_t)$$

* ↓
Minimize regret

\hat{y}_t = Combined
experts prediction

Loss made by the
best expert

* Halving algorithm: (Assuming one of the experts is always right)

- 1) • All experts have weight of 1
- Use majority vote as combined prediction \hat{y}_t
- If incur loss with label y_t and remove expert if prediction y_t is incorrect.
- For each expert where prediction doesn't match y_t , make weight zero.

$$\Rightarrow C_1 = n \quad (\text{Number of experts})$$

if $\hat{y}_t \neq y_t$

$$\text{then } C_{t+1} \leq \frac{1}{2} C_t$$

$$\frac{C_{t+1}}{C_1} \leq \left(\frac{1}{2}\right)^k \rightarrow \begin{matrix} \text{Number of mistakes} \\ \text{in } T \text{ trials} \end{matrix}$$

$$\frac{1}{2}^k \leq n$$

$$\text{Taking log on both sides, } k \leq \log_2 n \quad \left\{ \begin{matrix} \text{Upper bound of} \\ \text{number of} \\ \text{mistakes.} \end{matrix} \right\}$$

$$\Rightarrow R_T = \sum_{t=1}^T L_{0-1}(\text{Halving}(x_t), y_t)$$

$$\therefore R_T = k$$

$$\Rightarrow R_T \leq \log_2 n$$

?
One
expert
is always
right

(Algoin slides)

(No expert needs to be always right)

Weighted majority algorithm:

- Experts making mistakes have their weights decreased. Weights are all 1's initially.

$$w_i^{t+1} = w_i^t \cdot \exp(-n L_{0-1}(\varepsilon_i^t, y_t))$$

↑
Perform updates if $\hat{y}_t \neq y_t$ $n > 0$ is a parameter.
(This is to ensure small errors by experts are ignored since no expert is always right)

Theorem: m_T^* = Number of mistakes by best expert in T trials.

$$\text{Then } m_T^* \leq \frac{n m_T^* + \ln(n)}{\ln\left(\frac{2}{1+e^{-n}}\right)}$$

Mistakes made by weighted majority in T trials.

Proof:

Let $\hat{y}_t \neq y_t$

$$\begin{aligned} \text{then } w^{t+1} &= \sum_{j=1}^n w_j^{t+1} \\ &= \sum_{j=1}^n w_j^t \cdot e^{-n L_{0-1}(\varepsilon_j^t, y_t)} \end{aligned}$$

$$= \sum_{j: \varepsilon_j^t = y_t} w_j^t + \sum_{j: \varepsilon_j^t \neq y_t} w_j^t \cdot e^{-n}$$

($\sum w_j^t$ is the majority)

$$\because \hat{y}_t \neq y_t \quad = W_{\text{minority}} e^{-n} \cdot W_{\text{major}}$$

$$W^{t+1} \leq W_{\text{minority}} + e^{-n} W_{\text{major}}$$

We can upper bound it.

$$+ (W_{\text{major}} - W_{\text{minor}}) \cdot (1 - e^{-n})$$

$$\Rightarrow W^{t+1} \leq \frac{1 + e^{-n}}{2} \cdot W^t + \frac{(W_{\text{major}} - W_{\text{minor}})}{2} \cdot (1 - e^{-n})$$

the term.

$$\Rightarrow \frac{w^{t+1}}{w^t} < \frac{1+e^{-n}}{2}$$

$$\Rightarrow \frac{w^{T+1}}{w^1} = \prod_{t=1}^T \frac{w^{t+1}}{w^t}, \quad \text{Let } m_T \text{ be mistake in } T \text{ trials then}$$

$$\frac{w^{T+1}}{w^1} < \left(\frac{1+e^{-n}}{2} \right)^{m_T}$$

$$\therefore w^{T+1} < w^1 \cdot \left(\frac{1+e^{-n}}{2} \right)^{m_T} \rightarrow (1)$$

$$\Rightarrow \text{Now, } w^{T+1} = \sum_{i=1}^n w_i^{T+1} \geq \max_{i \in \{1, \dots, n\}} w_i$$

$$\Rightarrow w^{T+1} \geq \max_{i \in \{1, \dots, n\}} (w_i^1 \cdot e^{-n \sum_{t=1}^T L_{i,t}(\hat{E}_i^t)})$$

• This is a correct bound even with only k updates

$$w^{T+1} \geq e^{-n \min_{i \in \{1, \dots, n\}} \sum_{t=1}^T L_{i,t}(\hat{E}_i^t)}$$

m_T
(Number of mistakes made by best expert)

$$w^{T+1} \geq e^{-n \cdot m_T^*} \rightarrow (2)$$

- From (1), (2)

$$\Rightarrow e^{-n \cdot m_T^*} \leq n \left(\frac{1+e^{-n}}{2} \right)^{m_T}$$

$$\text{Taking log on both sides} \Rightarrow -n m_T^* \leq \ln(n) + m_T \ln\left(\frac{1+e^{-n}}{2}\right)$$

$$\Rightarrow m_T \ln\left(\frac{2}{1+e^{-n}}\right) \leq \ln(n) + n m_T^*$$

$$\Rightarrow m_T \leq \frac{\ln(n) + n m_T^*$$

$$\ln\left(\frac{2}{1+e^{-n}}\right)$$

~~⑤ Stochastic Algorithm~~

* 3) Exponential Weighted Average Algorithm:

- Assume prediction \hat{y}_t is from a convex set.
- Use a generalized loss: $\frac{|\hat{y}_t - y_t|}{2}$

• Prediction using, $\hat{y}_t = \sum_{i=1}^n p_i^t \cdot \varepsilon_i^t$

where $p_i^t = \frac{w_i^t}{\sum_{i=1}^n w_i^t}$

(Update every iteration)
• Weight update $w_i^{t+1} = w_i^t \cdot \exp(-\eta L(\varepsilon_i^t, y_t))$

$$\Rightarrow R_T = \sum_{t=1}^T L(\hat{y}_t, y_t) - \min_{i \in \{1, \dots, n\}} \sum_{t=1}^T L(y_t, \varepsilon_i^t)$$

Theorem: $\sum_{t=1}^T L(\hat{y}_t, y_t) \leq \min_{S \subseteq \{1, \dots, n\}} \sum_{i \in S} \sum_{t=1}^T L(\varepsilon_i^t, y_t) + \frac{\ln n}{n} + \frac{Tn}{8}$

$$R_T \leq Tn/8 + \frac{\ln n}{n}$$

(Proof on the next page) \rightarrow

Convex set: Let $X \subseteq \mathbb{R}^d$, X is convex if $\forall x_1, x_2 \in X$,
 $\lambda x_1 + (1-\lambda)x_2 \in X$ s.t. $0 \leq \lambda \leq 1$

Proof: $W_{t+1} = \sum_{i=1}^n w_i^{t+1}$

$$* \Rightarrow \frac{W_{t+1}}{W_t} = \frac{\sum_{i=1}^n w_i^t \cdot e^{-nL(\varepsilon_i^t, y_t)}}{\sum_{i=1}^n w_i^t}$$

$$= \sum_{i=1}^n p_i^t \cdot e^{-nL(\varepsilon_i^t, y_t)}$$

$$= \mathbb{E}_I [e^{-nL(\varepsilon_I, y_I)}]$$

• We know $\mathbb{E}[e^{\alpha x}] \leq e^{\alpha M} \cdot e^{\alpha^2/8}$

? (Hoeffding Lemma) $M = \mathbb{E}[x]$

$$0 \leq x \leq 1$$

Also $\frac{W_{t+1}}{W_t} \Rightarrow$ Our $L(\varepsilon_I, y_I)$ is bounded

between 0, 1 so we can apply Hoeffding

$$\frac{W_{t+1}}{W_t} \leq e^{-n\mathbb{E}_I(L(\varepsilon_I, y_I)) \cdot n^2/8}$$

~~Now, $\mathbb{E}_I[L(\varepsilon_I, y_I)] = \sum_{i=1}^n$~~

$$\prod_{t=1}^T \frac{W_{t+1}}{W_t} \leq e^{-n \cdot \sum_{t=1}^T \mathbb{E}_{I_t} [L(\varepsilon_{I_t}, y_t)] \cdot n^2/8}$$

• Log on both sides,

$$\log W_{T+1} - \log W_1 \leq \frac{Tn^2}{8} - n \sum_{t=1}^T \mathbb{E}_{I_t}$$

$$\because W_1 = n \quad , \quad \log W_{T+1} - \log n \leq \frac{Tn^2}{8} - n \sum_{t=1}^T \mathbb{E}_{I_t}$$

• From Jensen's Inequality, (for convex functions)

$$f(\mathbb{E}[z]) \leq \mathbb{E}[f(z)]$$

$$\therefore f(\sum p_i z_i) \leq \sum p_i f(z_i)$$

\Rightarrow Going back,

$$\log W_{T+1} - \log h \leq \frac{Tn^2}{8} - n \sum_{t=1}^T L(\mathbb{E}_{\mathcal{I}_t}[\varepsilon_{\mathcal{I}_t}], y_t)$$

$$\rightarrow \mathbb{E}_{\mathcal{I}_t}[\varepsilon_{\mathcal{I}_t}] = \sum_{i=1}^n p_i z_i^t - \bar{\varepsilon}_i^t = \hat{y}_t$$

$$\Rightarrow \log W_{T+1} - \log h \leq \frac{Tn^2}{8} - n \sum_{t=1}^T L(\hat{y}_t, y_t) \quad \rightarrow ①$$

$$\text{Now, } \Rightarrow \log W_{T+1} = \log \sum_{i=1}^n w_i^{T+1} \geq \max_{i \in \{1, \dots, n\}} \log w_i^{T+1}$$

$$\log W_{T+1} \geq \max_{i \in \{1, \dots, n\}} \log w_i^T \cdot e^{-n \sum_{t=1}^T L(\varepsilon_i^t, y_t)}$$

$$= \max_{i \in \{1, \dots, n\}} \log(w_i^T \cdot e^{-n \sum_{t=1}^T L(\varepsilon_i^t, y_t)})$$

$$= \log e^{-n \min_{i \in \{1, \dots, n\}} \sum_{t=1}^T L(\varepsilon_i^t, y_t)}$$

$$= -n \min_{i \in \{1, \dots, n\}} \sum_{t=1}^T L(\varepsilon_i^t, y_t) \quad \rightarrow ②$$

• From ①, ②,

$$-n \min_{i \in \{1, \dots, n\}} \sum_{t=1}^T L(\varepsilon_i^t, y_t) - \log h \leq \frac{Tn^2}{8}$$

$$-n \sum_{i \in \{1, \dots, n\}} \min_{t=1}^T L(\varepsilon_i^t, y_t) - n \sum_{i=1}^T \hat{y}_i$$

(Note, \hat{y}_i)

$$\Rightarrow R_T \leq \frac{Tn^2}{8} + \frac{\log h}{n}$$

(S. 8 total probability)

Note: Doubling trick to find optimal n in online learning.
 * (To minimize regret bound)

Note: Winnow Algorithm:

- Based on weighted majority
- Each feature can be thought of as an expert
- Positive weight for each feature
- Useful when there are a few distinguishable features but we don't know which.

⇒ Initial weights are $1/N$ per expert.

$$\hat{y}_i = \text{sign}[(w^\pm)^T \cdot x_i]$$

$$y_i \in \{+1, -1\}$$

$$\Rightarrow \text{Loss} = L_{0-1}(\hat{y}_i, y_i)$$

• If $\hat{y}_i \neq y_i$, update,

$$w_i^{\pm+1} = w_i^\pm \cdot \frac{\exp(n y_i x_i^\pm)}{Z_i}$$

$$\text{where } Z_i = \sum_{j=1}^n w_j^\pm \cdot \exp(n y_j x_j^\pm)$$

(6)

Theorem:

$$y_i \in \{+1, -1\}$$

$$x_i \in \mathbb{R}^N$$

$$R_\infty = \max_{i \in \{1, \dots, N\}} \|x_i^\pm\|_\infty \xrightarrow{\text{Samp Norm}}$$

$$\|x_i^\pm\|_\infty = \max_{i \in \{1, \dots, N\}} |x_{ii}^\pm| \xrightarrow{\text{Ideal weight}}$$

$$\|u\|_1 = \sum_{i=1}^n u_i \quad (\text{since all } u_i \geq 0)$$

• Assume γ such that $y_i \cdot u^\pm x_i^\pm \geq \gamma$

• $u \in \mathbb{R}^N, u_i \geq 0 \forall i \in \{1, \dots, N\}$

$$\|u\|_1 \leq N$$

$$\Rightarrow k \leq \frac{\|u\|_1 \ln N}{n\gamma - \|u\|_1 \ln \left(\frac{\exp(nR_\infty) + \exp(nR_\infty)}{2} \right)}$$

Mistake bound

(Optimizing wrt n)

$$k \leq \frac{2R_\infty^2 \|u\|_1^2}{\sqrt{2}} \ln N$$

$$\text{Ans} \Rightarrow \bar{P} = \frac{u}{\|u\|_1} \quad \left\{ \begin{array}{l} \text{The weight we} \\ \text{want to match (Ideal} \\ \text{weight)} \end{array} \right.$$

* We want to minimize KL divergence,

$$\text{Consider, } \text{KL}(\bar{P}, \bar{W}^t) - \text{KL}(\bar{P}, \bar{W}^{t+1})$$

$$= \sum_{i=1}^N p_i \ln \frac{p_i}{w_i^t} - \sum_{i=1}^N p_i \ln \frac{p_i}{w_i^{t+1}}$$

$$= \sum_{i=1}^N p_i \ln \frac{w_i^{t+1}}{w_i^t} = \sum_{i=1}^N p_i \ln \frac{e^{n y_t x_i^t}}{z_t}$$

$$= \sum_{i=1}^N p_i [n(y_t x_i^t - \ln z_t)]$$

$$= n y_t \sum_{i=1}^N p_i x_i^t - \ln z_t$$

$$= n y_t \bar{P} \cdot \bar{x}_t - \ln z_t$$

$$= n y_t \frac{u}{\|u\|_1} \cdot x_t - \ln z_t$$

From our assumption, $y_t \cdot u^T \cdot x^t \geq \gamma$

$$\therefore \text{KL}(\bar{P}, \bar{W}^t) - \text{KL}(\bar{P}, \bar{W}^{t+1})$$

$$\geq \frac{n \gamma}{\|u\|_1} - \ln z_t \rightarrow ①$$

$$\Rightarrow \text{Now, } z_t = \sum_{i=1}^N w_i^t e^{n y_t x_i^t}$$

$$= \sum_{i=1}^N w_i^t \cdot e^{\left[\frac{(1+n y_t x_i^t)/R_{00}}{2} \right] R_{00} + (-R)}$$

$$\text{And } \geq u \cdot u + \gamma u \cdot u$$

$$\Rightarrow z_t = \sum_{i=1}^N w_i^t \cdot e^{n \left[\frac{(1+(y_t x_i^t)/R_{00})}{2} R_{00} + (-R) \right]}$$

From, $f(\lambda z_1 + (1-\lambda) z_2) \leq \lambda f(z_1) + (1-\lambda) f(z_2)$ we can say

$$z^* \leq \sum_{i=1}^N w_i^* \left[\frac{1 + \frac{y_i x_i^*}{R_{\text{loss}}}}{2} \cdot e^{n R_{\text{loss}}} \right]$$

$$\Rightarrow z^* \leq \sum_{i=1}^N w_i^* \left[\frac{1 - \frac{y_i x_i^*}{R_{\text{loss}}}}{2} \cdot e^{-n R_{\text{loss}}} \right]$$

$$z^* \leq \sum_{i=1}^N w_i^* \left(\frac{e^{n R_{\text{loss}}} + e^{-n R_{\text{loss}}}}{2} \right) + \left(\frac{e^{n R_{\text{loss}}} - e^{-n R_{\text{loss}}}}{2} \right)$$

• $y_i \sum_{i=1}^N w_i^*$
-ve when misclassified

$$\therefore z^* \leq \sum_{i=1}^N w_i^* \left(\frac{e^{n R_{\text{loss}}} + e^{-n R_{\text{loss}}}}{2} \right)$$

$$z^* \leq \frac{e^{n R_{\text{loss}}} + e^{-n R_{\text{loss}}}}{2} \quad (\because \text{Weights sum up to 1}) \rightarrow (2)$$

From (1) (2)

$$\therefore KL(\bar{P}, \bar{W}^*) - KL(\bar{P}, \bar{W}^{T+1}) \geq \frac{n\gamma}{\|w\|_1} - \ln\left(\frac{e + e^{-n R_{\text{loss}}}}{2}\right) \rightarrow (3)$$

- If the algo makes k mistakes in T trials

$$KL(\bar{P}, \bar{W}^*) - KL(\bar{P}, \bar{W}^{T+1}) \geq k \left(\frac{n\gamma}{\|w\|_1} - \ln\left(\frac{e + e^{-n R_{\text{loss}}}}{2}\right) \right)$$

$$\text{Now, } KL(\bar{P}, \bar{W}^*) = \sum_{i=1}^N p_i \ln p_i - \sum p_i \ln \frac{1}{N}$$

$$= \underbrace{\sum p_i \ln p_i}_{-\text{ve}} + \ln N \leq \ln N$$

$$\therefore \ln N \geq k \left(\frac{n\gamma}{\|w\|_1} - \ln\left(\frac{e^{n R_{\text{loss}}} + e^{-n R_{\text{loss}}}}{2}\right) \right)$$

Now, differentiate w.r.t. γ to get γ^*

$$\Rightarrow \gamma^* = \frac{1}{2R_\infty} \ln \left(\frac{\|u\|_1 + \gamma}{\|u\|_1 - \gamma} \right)$$

$$\Rightarrow k \leq \frac{\ln N}{g\left(\frac{\gamma}{\|u\|_1}\right)} \rightarrow ①$$

$$\text{where } g(\epsilon) = \frac{(1+\epsilon)}{2} \ln(1+\epsilon)$$

$$\begin{aligned} \text{Should } \epsilon \leq 1? \\ + \frac{(1-\epsilon)}{2} \ln(1-\epsilon) \geq \frac{\epsilon^2}{2} \end{aligned}$$

Note: Cauchy-Schwarz: $a^T b \leq \|a\|_2 \|b\|_2$
extended $a^T b \leq \|a\|_1 \cdot \|b\|_\infty$

$$\begin{aligned} \text{Now, } \gamma &\leq \|u\|_1 \cdot \|u+x\|_\infty \\ &\leq \|u\|_1 \cdot R_\infty \end{aligned}$$

$$\Rightarrow \frac{\gamma}{\|u\|_1 \cdot R_\infty} \leq 1 \quad \left\{ \begin{array}{l} \text{This can be} \\ \text{put in } g(\epsilon) \end{array} \right.$$

$$\therefore \text{From } ①, \quad k \leq \frac{2 \ln N}{\frac{\gamma^2}{R_\infty^2 \cdot \|u\|_1^2}}$$

$$\Rightarrow k \leq \frac{2 R_\infty^2 \|u\|_1^2}{\gamma^2} - \ln N$$

$$\begin{aligned} \|u+x\|_\infty &\leq \|u\|_1 + \|x\|_\infty \\ \|u+x\|_\infty &\leq \|u\|_1 + R_\infty \end{aligned}$$

$$\|u+x\|_\infty \leq \|u\|_1 + R_\infty$$

* Online Gradient Descent: (General)

- General gradient descent methods where we use a sample at a time for inputs.

$$\Rightarrow \tilde{x}^{t+1} = x_t - \eta \nabla C_t(x_t) \quad \begin{matrix} \text{cost function can change} \\ \text{every} \\ \text{epoch} \end{matrix}$$

$\tilde{x}^{t+1} = P_{\Omega}(\tilde{x}^{t+1})$

where $P_{\Omega}(\tilde{x}^{t+1}) = \min_{x \in \Omega} \|x - \tilde{x}^{t+1}\|_2^2$

(Ω is our search space, so if \tilde{x}^{t+1} is outside search space, we adjust it accordingly)

Here, $R_T = \sum_{t=1}^T C_t(x_t) - \min_{x \in \Omega} \sum_{t=1}^T C_t(x)$

$$\Rightarrow R_T \leq \frac{1}{2} \left(\frac{D^2}{n} + nG^2 T \right) \quad \begin{matrix} \text{Best possible } x \\ \text{in entire} \\ \text{search space} \end{matrix}$$

Slides for details

Setting $\eta^* = \frac{D}{G\sqrt{T}}$

$$R_T \leq \frac{1}{2} D G \sqrt{T}$$

$$\Rightarrow \text{Since } C_t \text{ is convex, } C_t(x) - C_t(x_t) \geq \nabla C_t(x_t)^T (x - x_t)$$

$$\Rightarrow C_t(x) - C_t(x_t) \leq \nabla C_t(x_t)^T (x_t - x)$$

$$? \quad = \frac{x^t - \tilde{x}^{t+1}}{\eta} (x_t - x)$$

$$= \frac{1}{2n} \left[\|x - x^t\|^2 - \|x - \tilde{x}^{t+1}\|^2 + \|x^t - \tilde{x}^{t+1}\|^2 \right]$$

Now, $\|x - \tilde{x}^{t+1}\|_2 \geq \|x - x^t\|_2$
 $+ x \in \Omega$

$$C_{\pm}(x_{\pm}) - C_{\pm}(x) \leq \frac{1}{2n} \left(\|x - x^{\pm}\|^2 + \|x - x^{\pm+1}\|^2 \right)$$

$$= \frac{1}{2n} \left[\|x - x^{\pm}\|^2 - \|x - x^{\pm+1}\|^2 + n^2 \| \nabla C_{\pm}(x_{\pm}) \|^2 \right]$$

$$\therefore \sum_{\pm=1}^{\pm} C_{\pm}(x_{\pm}) - \sum_{\pm=1}^{\pm} C_{\pm}(x) \leq \frac{1}{2n} \left[\|x - x^{\pm}\|^2 - \|x - x^{\pm+1}\|^2 + n^2 \| \nabla C_{\pm}(x_{\pm}) \|^2 \right]$$

$$(A) + (B) + (C) \Rightarrow (A) + 2 \log T \leq \frac{1}{2n} \left[\|x - x^{\pm}\|^2 - \|x - x^{\pm+1}\|^2 + n^2 G^2 T \right]$$

$$(G = \| \nabla C_{\pm}(x_{\pm}) \|) \leq \frac{1}{2n} [D^2 + n^2 G^2 T]$$

De learned norm { $\| \cdot \|_p \times \| \cdot \|_q = \text{constant}$ }

$$\text{Note: } a^T b \leq \|a\|_p \|b\|_q \Rightarrow \frac{1}{p} + \frac{1}{q} = 1$$

$$\text{ex: } a^T b \leq \|a\|_1 \|b\|_\infty$$

we obtain at least one large value

large value

work 2001

already
done

- * Multi-Armed Bandits
- n arms, t denotes trial
 $\Rightarrow x_i^t$ is the reward obtained by pulling arm i on trial t

$x_i^t \in [0, 1]$ } Bounded, $x_i^t \sim q_{V_i}$ } Drawn from distribution q_{V_i}
 Arm pulled at t^{th} place

\Rightarrow Given $(i_1, x_{i_1}^1), (i_2, x_{i_2}^2), \dots, (i_{t-1}, x_{i_{t-1}}^{t-1})$ we need to select arm $i_t \in \{1, \dots, n\}$ to play.

* Regret $R_T(A)$ = $\max_{i \in \{1, \dots, n\}} G_T(i) - G_T(A)$

\Downarrow Expected regret

$$= \mathbb{E}_{(x^1, \dots, x^T)} [R_T(A)] \quad G_T(i) = \sum_{t=1}^T x_i^t$$

Arm selected by MAB

\Rightarrow Let $M_i = \mathbb{E}_{x_i \sim q_{V_i}} [x_i]$ Mean reward of arm i

$$M^* = \max_{i \in \{1, \dots, n\}} M_i, \Delta_i = M^* - M_i$$

- Expected regret is hard to handle, so we use pseudo regret

* $\tilde{R}_T(A) = \max_{i \in \{1, \dots, n\}} \mathbb{E}_{(x^1, \dots, x^T)} \left[\sum_{t=1}^T x_i^t - \sum_{t=1}^T x_{i_t}^t \right]$

Note: Show pseudo regret is upper bounded by expected regret

$(i_t = \text{Arm chosen by MAB algo at time } t)$

$$\tilde{R}_T(A) = T M^* - \sum_{t=1}^T \mathbb{E}_{(x^1, \dots, x^{t-1})} \mathbb{E}_{(x^t | x^1, \dots, x^{t-1})} [x_{i_t}^t]$$

$$\Rightarrow \tilde{R}_T(A) = T M^* - \sum_{t=1}^T \mathbb{E}_{(x^1, \dots, x^{t-1})} [x_t^A] \cdot \mathbb{E}_{x^t} [x_t^A]$$

(\because Arms draw rewards independent of past here)

$$= T M^* - \sum_{t=1}^T \mathbb{E}_{(x^1, \dots, x^{t-1})} [\mu_{i_t}]$$

$$\tilde{R}_T(A) = T M^* - \sum_{t=1}^T \mathbb{E} [\mu_{i_t}]$$

* Action Values: $\hat{\mu}_i^t$ = Estimated mean for arm i at epoch t
 $(\hat{\mu}_i^t = \frac{\sum_{s=1}^t \alpha_i^s \cdot \mathbb{I}_{\{i_s=i\}}}{\sum_{s=1}^t \mathbb{I}_{\{i_s=i\}}})$
 (known as action values)

\Rightarrow A simple greedy way of arm selection is select $i_t = \arg \max_{i \in \{1, \dots, n\}} \hat{\mu}_i^{t-1}$

\hookrightarrow This will never perform exploration though.

- This $\hat{\mu}_i^t$ can be derived just based on previous estimate (proof in slides)

$$? * \hat{\mu}_i^{t+1} = \hat{\mu}_i^t$$

* 1) E-Greedy Approach:

- In each trial, select best arm with prob $1-\epsilon$ & randomly select with prob ϵ

(Tradeoff b/w exploration & exploitation)

* 2) Upper confidence bound (UCB)

- At each interval, choose arm with longest UCB

- UCB of an arm can be high if its expected mean is high or it hasn't been sampled enough recently.

Note: $X \sim [0, 1]$, $\mu = \mathbb{E}(X) \Rightarrow (\mathbb{E}_{x_i \sim \mu} [x_i] = \mu_i)$

$$*\hat{\mu}_N = \frac{\sum_{i=1}^N x_i}{N}$$

From
Hoeffding Lemma

$$-2N\epsilon^2$$

$$\Rightarrow P[|\mu - \hat{\mu}_N| > \epsilon] \leq 2e^{-2N\epsilon^2}$$

* (As number of trials increase \rightarrow expected average goes closer to true average)

(Proof next page)

- Initially bounds are ~~large~~, since we don't have a lot of info but as we perform more trials it gets better.

$$\Rightarrow \text{Select arm } i^* \in \arg \max_{i \in \{1, \dots, n\}} \left[\hat{\mu}_i^{t+1} + \frac{\sqrt{2\ln t}}{2N_i^{t+1}} \right]$$

* We will prove why we use this later

⑧

N_i^t = Number of times arm i is pulled in trials.

$$\Rightarrow P[\mu_i - \hat{\mu}_i^t > \epsilon]$$

$$\Rightarrow P\left[\frac{1}{N_i^t} \sum_{k=1}^{N_i^t} x_{ik} - \frac{1}{N_i^t} \sum_{k=1}^{N_i^t} \hat{x}_{ik}^t > \epsilon\right]$$

Note: $\mathbb{E}\left[\frac{1}{N_i^t} \sum_{k=1}^{N_i^t} x_{ik}\right] = \frac{1}{N_i^t} \sum_{k=1}^{N_i^t} \mathbb{E}[x_{ik}] = \frac{1}{N_i^t} \sum_{k=1}^{N_i^t} \mu_i = \mu_i$

Assuming arm is pulled always pulled

else use indicator

$$\Rightarrow \sum_{s=1}^{N_i^t} \mathbb{I}_{\{x_{is} = \mu_i\}} = N_i^{t+1}$$

Consider,

$$P[\mu_i - \hat{\mu}_i^t > \epsilon] = P\left[\frac{\mathbb{E}\left[\sum_{s=1}^{N_i^t} x_{is} \mathbb{I}_{\{x_{is} = \mu_i\}}\right]}{N_i^t} - \frac{\sum_{s=1}^{N_i^t} \hat{x}_{is}^t \mathbb{I}_{\{x_{is} = \mu_i\}}}{N_i^t} > \epsilon\right]$$

$$= P\left[e^{\lambda \left\{ \mathbb{E}\left[\sum_{s=1}^{N_i^t} x_{is}^s \mathbb{I}_{\{x_{is} = \mu_i\}}\right] - \sum_{s=1}^{N_i^t} \hat{x}_{is}^s \mathbb{I}_{\{x_{is} = \mu_i\}} \right\}} - e^{\lambda N_i^t \epsilon} > \epsilon\right]$$

$$= P\left[e^{\lambda \left\{ \mathbb{E}\left[\sum_{s=1}^{N_i^t} x_{is}^s \mathbb{I}_{\{x_{is} = \mu_i\}}\right] - \sum_{s=1}^{N_i^t} \hat{x}_{is}^s \mathbb{I}_{\{x_{is} = \mu_i\}} \right\}} - e^{\lambda N_i^t \epsilon} > \epsilon\right]$$

We can do this because $P[X > \epsilon] \geq P[e^{X/\epsilon} > e]$

Since e^x is a monotonically increasing function.

$$= P\left[e^{\sum_{i=1}^{t-1} \{E[x_i] \mathbb{I}_{\{x_i=\bar{x}_i\}} - x_i^* \mathbb{I}_{\{x_i=\bar{x}_i\}}\}} > e^{dN_i^* \epsilon}\right]$$

$$\leq \frac{E\left[e^{\sum_{i=1}^{t-1} \{E[x_i] \mathbb{I}_{\{x_i=\bar{x}_i\}} - x_i^* \mathbb{I}_{\{x_i=\bar{x}_i\}}\}}\right]}{e^{dN_i^* \epsilon}}$$

Applying Markov's inequality

$$* P[X \geq a] \leq \frac{E[X]}{a}$$

Substituting \bar{x}_i in place of x_i in the above equation

$$\frac{E[\bar{x}_i]}{a} + b \leq \frac{E[\bar{x}_i]}{a} + \frac{b}{a}$$

$$\frac{E[\bar{x}_i]}{a} + b = \frac{E[\bar{x}_i]}{a} + \frac{b}{a} + \frac{b-a}{a}$$

$$\frac{E[\bar{x}_i]}{a} + b = \frac{E[\bar{x}_i]}{a} + \frac{b-a}{a}$$

$$\frac{E[\bar{x}_i]}{a} + b \leq \frac{E[\bar{x}_i]}{a} + \frac{b-a}{a}$$

$$\frac{E[\bar{x}_i]}{a} + b = \frac{E[\bar{x}_i]}{a} + \frac{b-a}{a}$$

to the right side
canceling a
both sides
canceling a

Upon solving, we get $P[|\bar{x}_i - \hat{x}_i^*| > \epsilon] \leq e^{-2N_i^* \epsilon^2}$

As trials increase, expected average gets closer to true average.

Lemma: Let i^* denote optimal arm $M_{i^*} = M^*$ and fix any $i: D_i > 0$. If UCB selects arm i at trial t ($i^* \neq i$) then atleast one of the following holds,

$$i) \hat{M}_{i^*}^{t-1} \leq M^* - \sqrt{\frac{\alpha \ln t}{2N_{i^*}^{t-1}}}$$

$$ii) \hat{M}_{i^*}^{t-1} \geq M_i + \sqrt{\frac{\alpha \ln t}{2N_i^{t-1}}}$$

$$iii) N_i^{t-1} \leq \frac{2\alpha \ln T}{D_i^2} \quad (D_i = M^* - M_i)$$

Proof:

(By contradiction) i) is false, $M^* < \hat{M}_{i^*}^{t-1} + \sqrt{\frac{\alpha \ln t}{2N_{i^*}^{t-1}}}$



$$\Rightarrow N_i^{t-1} > \frac{2\alpha \ln T}{D_i^2} \quad \{ \because iii) \text{ is false}$$

$$D_i > \sqrt{\frac{2\alpha \ln T}{N_i^{t-1}}}$$

$$\Rightarrow \hat{M}_{i^*}^{t-1} + \sqrt{\frac{\alpha \ln t}{2N_{i^*}^{t-1}}} > M_i + \underbrace{D_i}_{M^*}$$

$$\Rightarrow \hat{M}_{i^*}^{t-1} + \sqrt{\frac{\alpha \ln t}{2N_{i^*}^{t-1}}} > M_i + \sqrt{\frac{2\alpha \ln T}{N_i^{t-1}}}$$

(\because ii) is false)

$$\geq M_i + \sqrt{\frac{2\alpha \ln T}{N_i^{t-1}}}$$

$$\geq \hat{M}_i^{t-1} + \sqrt{\frac{2\alpha \ln T}{N_i^{t-1}}}$$

$$- \sqrt{\frac{\alpha \ln t}{2N_i^{t-1}}}$$

- This implies that the algorithm pulls the optimal arm but our assumption was

that arm i is pulled, hence its a contradiction.

$$\geq \hat{M}_i^{t-1} + \sqrt{\frac{\alpha \ln t}{2N_i^t}}$$

* Now, consider the pseudo regret $\tilde{R}_T(A)$

$$\begin{aligned} \tilde{R}_T &= T\mu^* - \mathbb{E}\left[\sum_{t=1}^T \mu_{i_t}\right] \quad (\text{From 2 pages earlier}) \\ &= T\mu^* - \mathbb{E}\left[\sum_{i=1}^n N_i^T \mu_i\right] \end{aligned}$$

($N_i^T = \text{Number of times arm } i \text{ is pulled in } T \text{ trials}$)

$$= T\mu^* - \sum_{i=1}^n \mathbb{E}[N_i^T]$$

~~μ^*~~

$$= \sum_{i=1}^n \mathbb{E}[N_i^T] \cdot \mu^* - \sum_{i=1}^n \mathbb{E}[\mu_i N_i^T]$$

$$= \sum_{i=1}^n \mathbb{E}[N_i^T (\mu^* - \mu_i)]$$

$$= \sum_{i=1}^n \mathbb{E}[N_i^T \Delta_i] \quad \begin{array}{l} \text{Fixed value} \\ \text{based on actual means} \end{array}$$

$$= \sum_{i=1}^n \Delta_i \mathbb{E}[N_i^T]$$

⑨

$$\Pr[|\mu - \hat{\mu}_i| \geq t] \leq 2e^{-2N_i t^2}$$

$$\Pr[|\mu_i - \hat{\mu}_i^*| \geq \epsilon] \leq e^{-2N_i \epsilon^2} \quad \begin{array}{l} \text{Proof from} \\ \text{Hoeffding Lemma} \end{array}$$

$$\hat{\mu}_i = \frac{\sum x_{i,t}}{N} \quad \forall t \in [0, 1], \mu = \mathbb{E}[x]$$

$$\hat{\mu}_i^* = \frac{\sum_{s=1}^{t-1} x_{i,s} \mathbb{I}_{\{x_{i,s}=i\}}}{\sum_{s=1}^{t-1} \mathbb{I}_{\{x_{i,s}=i\}}} \rightarrow N_i^{t-1}$$

Proof: Complete proof in lecture notes (on Moodle)

* UCB Arm Selection:

- $P[\mu_i - \hat{\mu}_i^* \geq \epsilon] \leq e^{-2N_i^{t-1}\epsilon^2} \leq \delta_t$

$$-2N_i^{t-1}\epsilon^2 \leq \ln \delta_t$$

$$\ln \frac{1}{\delta_t} \leq 2N_i^{t-1}\epsilon^2$$

Assume
its an
approx.

$$\epsilon \geq \sqrt{\frac{\ln(1/\delta_t)}{2N_i^{t-1}}} \rightarrow \epsilon^*$$

$$\Rightarrow \text{Now } P[\mu_i - \hat{\mu}_i^* \geq \sqrt{\frac{\ln(1/\delta_t)}{2N_i^{t-1}}}] \leq \delta_t$$

$$\Rightarrow \text{Consider, } P[\mu_i - \hat{\mu}_i < \epsilon] > 1 - \delta_t$$

that is, $\mu_i - \hat{\mu}_i < \epsilon^*$ with prob
atleast $1 - \delta_t$

$\mu_i < \hat{\mu}_i + \epsilon^*$ with prob $1 - \delta_t$

$\mu_i < \hat{\mu}_i + \sqrt{\frac{\ln(1/\delta_t)}{2N_i^{t-1}}}$ with prob $1 - \delta_t$

- Choose $\delta_t = \frac{1}{t^\alpha}$ \rightarrow Just has to be a function of t

* $\Rightarrow \boxed{\mu_i < \hat{\mu}_i + \sqrt{\frac{\alpha \ln t}{2N_i^{t-1}}}}$ with prob $1 - \delta_t$
which is high.

\therefore We use this as the criterion for
arm selection in our UCB algo.

* UCB Regret Bound:

$$R_T = T\mu^* - \sum_{t=1}^T \mathbb{E}[\mu_{i_t}]$$

$$T = \mathbb{E}\left[\sum_{t=1}^T Z_{i_t}^T\right]$$

$$\Rightarrow R_T = \mathbb{E}\left[\sum_{t=1}^T Z_{i_t}^T\right] \mu^* - \sum_{t=1}^T \mathbb{E}[\mu_{i_t}]$$

$$\Rightarrow R_T = \mathbb{E}\left[\sum_{i=1}^n Z_{i_t}^T\right] \mu^* - \mathbb{E}\left[\sum_{t=1}^T \mu_{i_t}\right]$$

$$= \mathbb{E}\left[\sum_{i=1}^n Z_{i_t}^T\right] \mu^* - \mathbb{E}\left[\sum_{i=1}^n Z_{i_t}^T \mu_i\right]$$

$$\left(\because \sum_{t=1}^T \mu_{i_t} = \sum_{t=1}^T \sum_{i=1}^n \mu_i \mathbb{I}_{\{i_t=i\}} \right)$$

$$= \sum_{i=1}^n \mu_i \left(\sum_{t=1}^T \mathbb{I}_{\{i_t=i\}} \right) \xrightarrow{\text{let } N_i^T}$$

$$= \sum_{i=1}^n \mu_i N_i^T$$

$$\Rightarrow R_T = \mathbb{E}\left[\sum_{i=1}^n Z_{i_t}^T (\mu^* - \mu_i)\right]$$

$$= \mathbb{E}\left[\sum_{i=1}^n Z_{i_t}^T \Delta_i\right]$$

$$R_T = \sum_{i=1}^n \mathbb{E}[Z_{i_t}^T \Delta_i] \rightarrow \textcircled{1}$$

$$\text{Now, } \mathbb{E}[Z_{i_t}^T] = \mathbb{E}\left[\sum_{t=1}^T \mathbb{I}_{\{i_t=i\}}\right]$$

$$= \mathbb{E}\left[\sum_{t=1}^T \left(\underbrace{\mathbb{I}_{\{\sum_{i_t=i}, N_i^{t-1} \leq t_0\}}}_{\{i_t=i, N_i^{t-1} \leq t_0\}} + \mathbb{I}_{\{i_t=i, N_i^{t-1} > t_0\}} \right) \right]$$

↑ to (Upper bound)

$\because Z_i^t < \mu^*$
always.

$$\leq \sum_{t=t_0+1}^T \mathbb{E}\left[\sum_{i=1}^n \mathbb{I}_{\{i_t=i, N_i^{t-1} > t_0\}}\right]$$

$$\leq t_0 + \sum_{t=t_0+1}^T P[\hat{A}_{i^*}^t = i^*, N_{i^*}^{t-1} > t_0]$$

Choose $t_0 = \sqrt{\frac{2\alpha \ln T}{\Delta_{i^*}^2}}$, $\alpha = \text{Just a parameter.}$

$$\therefore E[N_{i^*}^T] \leq t_0 + \sum_{t=t_0+1}^T P[\text{event 1 or event 2}]$$

$$\leq t_0 + \sum_{t=t_0+1}^T [P[\text{event 1}] + P[\text{event 2}]]$$

$$\text{Now, } P[\text{event 1}] = P[\hat{A}_{i^*}^t \leq i^* - \sqrt{\frac{\alpha \ln t}{2N_{i^*}^{t-1}}}]$$

(Since the 3 condition lemma must hold)

$$P[\text{event 1}] = P\left[\bigcup_{s=1}^{t-1} \{\hat{A}_{i^*}^s \leq i^* - \sqrt{\frac{\alpha \ln s}{2s}}\}\right]$$

Assuming i^* was pulled any number of times from 1 to $t-1$.

$$P[\text{event 1}] \leq \sum_{s=1}^{t-1} P\{\hat{A}_{i^*}^s \leq i^* - \sqrt{\frac{\alpha \ln s}{2s}}\}$$

From Hoeffding's lemma, $P\{\hat{A}_{i^*}^s - M^* \leq$

$$-\sqrt{\frac{\alpha \ln s}{2s}}$$

$$= 1/t^\alpha$$

$$\therefore P[\text{event 1}] \leq \sum_{s=1}^{t-1} 1/t^{\alpha-1} \leq 1/t^{\alpha-1}$$

$$\leq \sum_{s=1}^t 1/t^{\alpha-1}$$

$$\therefore E[N_{i^*}^T] \leq t_0 + \frac{1}{\alpha-2} + \frac{1}{\alpha-2}$$

$$\therefore \left(\sum_{t=t_0+1}^T P[\text{event 1}] \leq \frac{1}{\alpha-2} \right)$$

\Rightarrow Similarly we can derive $\sum_{t=t_0+1}^T P[\text{event 2}] \leq \frac{1}{\alpha-2}$

$$\therefore \mathbb{E}[N_i T] \leq t_0 + \frac{2}{\alpha-2} = \frac{2\alpha \ln T}{\Delta_i^2} + 1 + \frac{2}{\alpha-2}$$

$$= \frac{2\alpha \ln T}{\Delta_i^2} + \frac{\alpha}{\alpha-2}$$

Now, going back to ①,

$$\tilde{R}_T = \sum_{i=1}^n \mathbb{E}[N_i T \Delta_i]$$

$$= \sum_{i=1}^n \Delta_i \mathbb{E}[N_i T]$$

$$\leq \sum_{i=1}^n \Delta_i \left(\frac{2\alpha \ln T}{\Delta_i^2} + \frac{\alpha}{\alpha-2} \right)$$

$$= \sum_{i: \Delta_i > 0} \left(\frac{2\alpha \ln T}{\Delta_i} \right) + \sum_{i: \Delta_i > 0} \frac{\alpha}{\alpha-2}$$

$$R_T[\text{UCB}(\alpha)] \leq \sum_{i: \Delta_i > 0} \left(\frac{2\alpha}{\Delta_i} \right) \ln T + \sum_{i: \Delta_i > 0} \frac{\alpha}{\alpha-2}$$

⑩ Adversarial MAB: (Exp 3)

- Need to choose arms non-deterministically.
- Actions selected randomly based on weights of each arm (CDF)

\hookrightarrow Let \tilde{l}_{i*}^+ be the loss corresponding to arm $i*$.

$$\Rightarrow \tilde{l}_{i*}^+ = \frac{\tilde{l}_{i*}^+}{p_{i*}} \prod_{j \neq i*} \{ j \neq i* \} \quad \begin{cases} \text{Importance} \\ \text{Sampling} \end{cases}$$

EXP3(Alg) - Adversarial MAB

$$p^t = \frac{1}{n} [1, \dots]$$

for $t = 1 \dots T$:

Observe $l_{it}^* \in [0, 1]$

for $i = 1 \dots n$:

$$\tilde{l}_i^t = \frac{l_i^*}{p_i^t} \mathbb{I}_{\{i=t\}}$$

$$\tilde{l}_i^t = \tilde{l}_i^{t-1} + \tilde{l}_i^t$$

for $i = 1 \dots n$:

$$p_i^{t+1} = \frac{e^{-n_t \tilde{l}_i^t}}{\sum_{j=1}^n e^{-n_t \tilde{l}_j^t}}$$

Note: $E_{l_* \sim p^t} [\tilde{l}_i^t] = E_{l_* \sim p^t} \left[\frac{l_i^*}{p_i^t} \mathbb{I}_{\{i=t\}} \right]$

*



$$= \sum_{j=1}^n p_j^t \cdot \frac{l_i^*}{p_i^t} \mathbb{I}_{\{j=i\}}$$

$$= l_i^*$$

\tilde{l}_i^t is an unbiased estimator of l_i^*

*

Regret Analysis:

$$\phi_t = -\frac{1}{n} \log \sum_{i=1}^n e^{-n_t \tilde{l}_i^t} \quad \text{Potential function}$$

$$\Rightarrow \phi_{t+1} - \phi_t = -\frac{1}{n} \log \left[\frac{\sum_{i=1}^n e^{-n_t \tilde{l}_i^{t+1}}}{\sum_{j=1}^n e^{-n_t \tilde{l}_j^t}} \right]$$

$$\Rightarrow \phi_{t+1} - \phi_t = -\frac{1}{n} \log \left[\frac{\sum_{i=1}^n e^{-n\tilde{l}_i^*} \cdot e^{-n\tilde{l}_i^*}}{\sum_{j=1}^n e^{-n\tilde{l}_j^*}} \right]$$

$$= -\frac{1}{n} \log \left[\sum_{i=1}^n e^{-n\tilde{l}_i^*} \cdot p_i^* \right]$$

$$= -\frac{1}{n} \log \left[\mathbb{E}_{i \sim p^*} [e^{-n\tilde{l}_i^*}] \right] \rightarrow 0$$

$$\text{Now, } e^{-x} \leq 1-x + \frac{x^2}{2} \quad \forall x \geq 0$$

(∴ The rest of the terms will end up summing to a -ve value)

Putting this in ①,

$$\phi_{t+1} - \phi_t \geq -\frac{1}{n} \log \left[\mathbb{E}_{i \sim p^*} \left[1 - \frac{n\tilde{l}_i^*}{2} + \frac{n^2(\tilde{l}_i^*)^2}{4} \right] \right]$$

$$\phi_{t+1} - \phi_t \geq -\frac{1}{n} \log \left[1 - \mathbb{E} \left[n\tilde{l}_i^* - \frac{n^2}{2} (\tilde{l}_i^*)^2 \right] \right]$$

$$\text{Now, } \log(1-x) \leq -x \Rightarrow -\log(1-x) \geq x$$

$$\therefore \phi_{t+1} - \phi_t \geq \frac{1}{n} \mathbb{E}_{i \sim p^*} \left[n\tilde{l}_i^* - \frac{n^2}{2} (\tilde{l}_i^*)^2 \right]$$

$$= \mathbb{E}_{i \sim p^*} [\tilde{l}_i^*] - \frac{n}{2} \mathbb{E}_{i \sim p^*} [(\tilde{l}_i^*)^2]$$

$$= \sum_{i=1}^n p_i^* \tilde{l}_i^* - \frac{n}{2} \sum_{i=1}^n p_i^* \cdot (\tilde{l}_i^*)^2$$

$$\phi_{t+1} - \phi_t \geq \sum_{i=1}^n l_i^+ \mathbb{I}_{\{j_t=i\}} - \frac{n}{2} \sum_{i=1}^n \frac{(l_i^+)^2}{p_i^+}$$

$$\Rightarrow \mathbb{E}[\phi_{t+1} - \phi_t | H_{t-1}] \geq \sum_{i=1}^n l_i^+ \mathbb{I}_{\{j_t=i\}}$$

~~$\cdot \mathbb{E}[l_i^+ | H_{t-1}]$~~

$$\mathbb{E}[\phi_{t+1} - \phi_t | H_{t-1}] \geq \mathbb{E}\left[\sum_{i=1}^n l_i^+ \mathbb{I}_{\{j_t=i\}}\right] - \frac{n}{2} \sum_{i=1}^n \frac{(l_i^+)^2}{p_i^+}$$

$$= \sum_{i=1}^n l_i^+ \cdot \mathbb{E}[\mathbb{I}_{\{j_t=i\}}]$$

$$- \frac{n}{2} \sum_{i=1}^n \frac{(l_i^+)^2}{p_i^+} \quad \rightarrow (2)$$

$$\text{Now, } \mathbb{E}_{j_t \sim p_i^+} [\mathbb{I}_{\{j_t=i\}}]$$

$$= \sum_{i=1}^n p_i^+ \mathbb{I}_{\{j_t=i\}} = p_i^+$$

Now, from (2),

$$\mathbb{E}[\phi_{t+1} - \phi_t | H_{t-1}] \geq \sum_{i=1}^n l_i^+ p_i^+ - \frac{n}{2} \sum_{i=1}^n (l_i^+)^2$$

$\because l_i^+ \in [0, 1] \quad \Rightarrow \quad \sum_{i=1}^n l_i^+ p_i^+ - \frac{n}{2} \cdot n$

$$\text{Now, } \mathbb{E}[\phi_{t+1} - \phi_t] \geq \mathbb{E}_{H_{t-1}} \left[\sum_{i=1}^n l_i^+ p_i^+ - \frac{n}{2} \cdot n \right]$$

$$\mathbb{E}[\phi_{T+1} - \phi_1] \geq \mathbb{E}[l_{i^*}^+] - \frac{n\bar{n}}{2}$$

Sum from 1 to T,

$$\mathbb{E}[\phi_{T+1} - \phi_1] \geq \mathbb{E}\left[\sum_{t=1}^T l_{i_t}^+\right] - \frac{nT\bar{n}}{2} \quad \rightarrow (3)$$

$$\text{Now, } \phi_{T+1} - \phi_1 = -\frac{1}{n} \log \sum_{i=1}^n e^{-n\tilde{L}_i^T}$$

$$\leq -\frac{1}{n} \log \left(\frac{e^{-n\tilde{L}_{i^*}^T}}{n} \right)$$

$$\text{where } i^* = \operatorname{argmin}_i \tilde{L}_i^T$$

$$\Rightarrow \phi_{T+1} - \phi_1 \leq \tilde{L}_{i^*}^T + \frac{1}{n} \log n$$

$$\mathbb{E}[\phi_{T+1} - \phi_1] \leq \tilde{L}_{i^*}^T + \frac{1}{n} \log n \quad \rightarrow (4)$$

Comparing LB & UB, (3) & (4)

$$\mathbb{E}\left[\sum_{t=1}^T l_{i_t}^+\right] - \frac{nT\bar{n}}{2} \leq \tilde{L}_{i^*}^T + \frac{\log n}{n}$$

Note: $\mathbb{E}[\tilde{L}_i^T] = \tilde{L}_i^T \quad (\text{For any } i)$

$$\therefore \mathbb{E}[\tilde{L}_i^T] = \mathbb{E}_{\{i_1, \dots, i_T\}} \left[\sum_{t=1}^T \frac{l_{i_t}^+}{P_{i_t}^+} \mathbb{I}_{\{i_t=i\}} \right]$$

$$= \sum_{t=1}^T \mathbb{E}_{\{i_t\}} \left[\frac{l_{i_t}^+}{P_{i_t}^+} \mathbb{I}_{\{i_t=i\}} \right]$$

$$= \sum_{t=1}^T l_{i^*}^+ = \tilde{L}_{i^*}^T$$

$$\Rightarrow \underbrace{\mathbb{E} \left[\sum_{t=1}^T l_i^* \right] - L_{i*}^T}_{\text{Regret}} \leq \frac{\log n}{n} + \frac{\eta \cdot T n}{2}$$

Differentiating w.r.t η & maximizing

$$\eta^* = \sqrt{\frac{2 \log n}{T n}}$$

$$\text{Regret with } \eta^* = \sqrt{2 T n \log n}$$

(In non adversarial, we get $\text{Regret}(\eta^*)$)

$$= \sqrt{2 T \log n}$$

\Rightarrow This is a partial information setting
so we have to pay the price of \sqrt{n} .

Note: EXT3 can be extended to contextual MAB setting. (Run EXT3 for each context available)

⑪

Reinforcement Learning:



- S = Set of States

- R = Set of Rewards

- $A(s)$ = Action set when state is s .

$\Rightarrow R_{t+1}$ = Reward at $t+1$ by action A_t being performed in state S_t .

- Return $G_{1:T} = R_{t+1} + \dots + R_{T \rightarrow \text{Time horizon}}$

\rightarrow Maximize this return.

⇒ Episodic tasks - Finite time horizon & reset and restart game every episode.

⇒ Continuous tasks - Infinite time horizon.

$$G_T = \sum_{k=0}^{\infty} \gamma^k \cdot R_{t+k-1}$$

Discounted return $\gamma \in [0, 1)$

* ($\Rightarrow \gamma = 0$, we focus only on immediate reward)
 (Myopic setting)

⇒ Markov ~~property~~ - Memoryless, only takes previous state & action into consideration

$$\begin{aligned} p_\pi [S_{t+1} = s', R_{t+1} = r | S_0, A_0, R_1, S_1, \dots, R_t, S_t, A_t] \\ = p_\pi [S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a] \end{aligned} \rightarrow ①$$

* Markov Decision Process (MDP):

↳ Reinforcement learning tasks which satisfy the markov property ① are called Markov decision processes.

Quantities of Interest → i) Expected Reward:

$$\begin{aligned} q_\pi(s, a) &= \mathbb{E}[R_{t+1} | S_t = s, A_t = a] \\ &= \sum_r q_r \cdot p_\pi[R_{t+1} = r | S_t = s, A_t = a] \end{aligned}$$

$$\begin{aligned} * \Rightarrow q_\pi(s, a) &= \sum_{r \in R} \sum_{s' \in S} p_\pi[R_{t+1} = r, S_{t+1} = s' | S_t = s, A_t = a] \\ &= \sum_{r \in R} \sum_{s' \in S} p[r, s' | s, a] \end{aligned}$$

i) State transition probability:

$$P[S_{t+1} = s' | S_t = s, A_t = a]$$

$$= \sum_{g \in R} P[S_{t+1} = g, R_{t+1} = r_t | S_t = s, A_t = a]$$

$$= \sum_{g \in R} P[g, r_t | s, a]$$

* ii) Expected reward (for state-action next state triplets):

$$\mathbb{E}[R_{t+1} | S_t = s, A_t = a]$$

ex: Recycling Robot

$$S = \{\text{high}, \text{low}\}$$

$$A(\text{high}) = \{\text{search}, \text{wait}\}$$

$$A(\text{low}) = \{\text{search}, \text{wait}, \text{recharge}\}$$

$S_t(s)$	$S_{t+1}(s)$	$A_t(a)$	$P(s' s, a)$	$r_t(s, a)$
high	high	search	α	r_{search}
high	low	search	$1 - \alpha$	r_{search}
high	high	wait	1	r_{wait}
high	low	wait	0	r_{wait}
high	high	search	$1 - \beta$	-3
high	low	search	β	r_{search}
Assume it recharge somehow	low	wait	0	r_{wait}
low	high	wait	1	r_{wait}
low	low	wait	0	0
low	high	recharge	0	0
low	low	recharge	1	0

* Policy: A policy π is a mapping from a set of states and a set of actions to probabilities

$$\pi(a|s) = P[A_t=a | S_t=s]$$

i) State-Value function:

$$V_\pi(s) = \mathbb{E}_\pi [G_t | S_t=s]$$

↓
Discounted
return

ii) State-Action pair value function:

$$* \cdot V_\pi(s, a) = \mathbb{E}_\pi [G_t | S_t=s, A_t=a]$$

$$V_\pi(s) = \mathbb{E}_\pi [G_t | S_t=s]$$

$$= \mathbb{E} [\mathbb{E}_\pi [G_t | S_t=s, A_t=a]]$$

$$= \sum_{a \in A(s)} \pi(a|s) V_\pi(s, a)$$

(2) $V_\pi(s) = \mathbb{E}_\pi [G_t | S_t=s]$

$$= \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t=s \right]$$

$$= \mathbb{E}_\pi \left[R_{t+1} + \sum_{k=1}^{\infty} \gamma^k R_{t+k+1} | S_t=s \right]$$

Also,

$$V_\pi(s) = \mathbb{E}_{A_t|S_t} [\mathbb{E}_\pi [G_t | S_t=s, A_t=a]]$$

$$= \mathbb{E}_{A_t|S_t} [V_\pi(s, a)]$$

$$= \sum_{a \in A(s)} V_\pi(s, a) \cdot P_\pi(A_t=a | S_t=s)$$

$$= \sum_{a \in A(s)} \pi(a|s) \cdot V_\pi(s, a)$$

$$\begin{aligned}
\Rightarrow q_{\pi}(s, a) &= \mathbb{E}_{\pi} [G_t | S_t = s, A_t = a] \\
&= \mathbb{E}_{\pi} \left[R_{t+1} + \sum_{k=1}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right] \\
&\quad \cancel{=} \cancel{\mathbb{E}_{S_{t+1}=s \mid S_t=s, A_t=a} \left[\mathbb{E}_{R_{t+1}+ \sum_{k=1}^{\infty} \gamma^k R_{t+k+1} \mid S_t=s, A_t=a} \right]} \\
&= \mathbb{E}_{\pi} [R_{t+1} \mid S_t = s, A_t = a] \\
&\quad + \mathbb{E}_{S_{t+1}=s' \mid S_t=s, A_t=a} \left[\mathbb{E}_{\pi} [\gamma G_{t+1} \mid S_t = s, A_t = a, S_{t+1} = s'] \right] \\
&= r(s, a) + \mathbb{E}_{S_{t+1}=s \mid S_t=s, A_t=a} \left[\mathbb{E}_{\pi} [\gamma G_{t+1} \mid S_{t+1}=s'] \right] \\
&= r(s, a) + \mathbb{E}_{S_{t+1}=s \mid S_t=s, A_t=a} \left[\gamma \cdot V_{\pi}(s') \right] \\
&= r(s, a) + \gamma \cdot \sum_{s' \in S} V_{\pi}(s') - p_{\pi} \left[S_{t+1}=s' \mid S_t=s, A_t=a \right] \\
&= r(s, a) + \gamma \sum_{s' \in S} V_{\pi}(s') \cdot p(s'|s, a) \\
&= r(s, a) + \gamma \sum_{s' \in S} V_{\pi}(s') \cdot \sum_{a'} p(a, s'|s, a)
\end{aligned}$$

* we have,

$$V_{\pi}(s) = \sum_{a \in A(s)} \pi(a|s) \left[r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V_{\pi}(s') \right]$$

* \downarrow

Bellman equation for state value function.

* Optimal Value function:

- The optimal state value function $V_\pi(s)$ is as follows,

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

Optimal State-Action value function:

$$q_\pi(s, a) = \max_{\pi} q_{\pi'}(s, a)$$

- Both of these maximums occur for the same policy π^*

$$V_{\pi^*}(s) = \max_{\pi} V_\pi(s) = \max_{\pi} \sum_{a \in A(s)} \pi(a|s) \cdot q_{\pi^*}(s, a)$$

Bellman optimality equation for state value

$$= \max_{a \in A(s)} q_{\pi^*}(s, a)$$

(∴ If we know q_π then $\pi(a|s) = 1$ for some

$a \in A(s)$ & $\pi(a_i|s) = 0$ for the rest)
 The best action to choose

$$q_{\pi^*}(s, a) = \max_{\pi} q_{\pi^*}(s, a) = \max_{\pi} q_{\pi^*}(s, a) + \sum_{s' \in S} p(s'|s, a) V_{\pi^*}(s')$$

$$= \max_{\pi} q_{\pi^*}(s, a) + \sum_{s' \in S} p(s'|s, a) - \sum_{a' \in A(s)} \pi(a'|s) q_{\pi^*}(s, a')$$

$$= \max_{a' \in A(s)} q_{\pi^*}(s, a') + \sum_{s' \in S} p(s'|s, a') \cdot q_{\pi^*}(s', a')$$

Bellman optimality equation for state-action value.

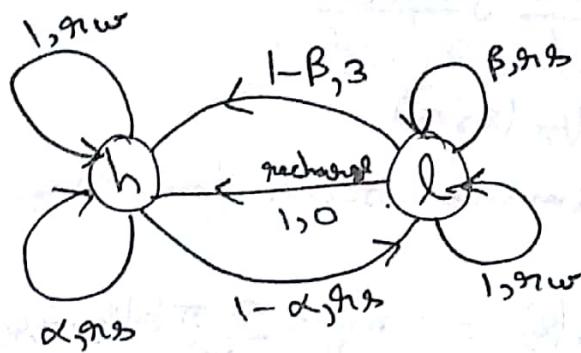
ex: Recycling robot (As we saw earlier)

$$S = \{l, h\}$$

$$A(l) = \{s, w, r\}$$

$$A(h) = \{s, w\}$$

s_t	s_{t+1}	a_t	$P(s' s,a)$	$r(s')$
h	h	s	α	r_s
h	l	s	$1-\alpha$	r_s
l	h	s	$1-\beta$	-3
l	l	s	β	r_s



(13)

$$V_\pi(s) = \max_{a \in A(s)} q_{\pi^*}(s, a)$$

$$F^\pi(s) = \max_{a \in A(s)} q_\pi(s, a) + \sum_{s'} p(s'|s, a) V_\pi(s')$$

~~$$V^\pi(s) = \max_a \{ p(h|s, a) \}$$~~

→ Bellman optimality operator

$$V^\pi = g^\pi + \gamma P^\pi \cdot V^\pi$$

Considering actions at each step are known.

Bellman's expectation operator

$$F^\pi(s) = g^\pi + P^\pi \cdot V$$

Set of all states

HW1: Note: Verify that Bellman Optimality operator satisfies monotonicity property.

HW2: Show that Bellman Optimality operator is γ -contraction mapping.

* Use: $\max_a f(a) - \max_a g(a) \leq \max_a (f(a) - g(a))$

* Properties of Bellman ops.: (expectation operator)

i) Monotonicity $\rightarrow V \leq U \Rightarrow F(V) \leq F(U)$

$$\begin{aligned} F(V) &= \gamma^\pi + \gamma P^\pi V \\ F(U) &= \gamma^\pi + \gamma P^\pi U \end{aligned}$$

$$\text{Difference: } F(V) - F(U) = \gamma P^\pi (V - U)$$

ii) Offset:

$$F(V + ce) = F(V) + \gamma ce$$

$$\leq \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Vector of all 1's

* iii) γ -contraction in L_∞ norm:

$$\Rightarrow \|F(V) - F(U)\|_\infty \leq \gamma \|V - U\|_\infty$$

$$\text{Lower bound: } F^\pi(V) - F^\pi(U) = \gamma P^\pi(V - U)$$

$$\text{Note: } \|V\|_\infty = \max_{i \in \{1, \dots, N\}} |V_i|$$

$$\Rightarrow \|F^\pi(V) - F^\pi(U)\|_\infty = \|\gamma P^\pi(V - U)\|_\infty$$

$$= \gamma \max_{j=1, \dots, N} \sum_i p_{ij}^\pi (V(s_i) - U(s_i))$$

$$\begin{aligned} \gamma P^\pi(V - U) &= \left[\sum_i p_{i,1}^\pi (V(s_1) - U(s_1)) \right. \\ &\quad \left. + \sum_i p_{i,2}^\pi (V(s_2) - U(s_2)) \right. \\ &\quad \left. + \cdots + \sum_i p_{i,N}^\pi (V(s_N) - U(s_N)) \right] \end{aligned}$$

$$\begin{aligned} &\leq \gamma \max_{i \in \{1, \dots, N\}} |V(s_i) - U(s_i)| \cdot (\sum_i p_{i,j}^\pi) \\ &\leq \gamma \max_i |V(s_i) - U(s_i)| \\ &\leq \gamma \|V - U\|_\infty \end{aligned}$$

So $\|F^\pi\|_\infty = \gamma \|P^\pi\|_\infty$

* Banach's Fixed Point theorem:

- Ensures the existence of unique fixed point (convex)

$$\hookrightarrow \lim_{k \rightarrow \infty} (F^\pi)^k(v) = v^\pi$$

$$\lim_{k \rightarrow \infty} (F)^k(v) = v^* \rightarrow \text{Solution to bellman's optimality equations.}$$

1) Policy Iteration: (Details in slides)

* For a given policy, we try to find the values of each state (estimated rewards)

\Rightarrow This gives us the best reachable reward from each state, we then update the policy

\hookrightarrow The policy is updated as such, (policy improvement)
 if $\exists a \in A(s) \text{ s.t. } q_\pi(s, a) \geq v_\pi(s)$
 then $\pi^1(s) = a$

$$\Rightarrow \pi_{k+1}(s) = \arg \max_{a \in A(s)} \left[q_\pi(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) \cdot v_\pi(s') \right]$$

(Steps are policy evaluation & policy improvement
 and then repeat)

\Rightarrow We stop at k^{th} iteration when $v_{\pi_{k+1}}(s) = v_{\pi_k}(s)$

• No of possible policies = $|A|^{|S|} \rightarrow$ No. of states.

↓
 Finite number
 and hence —

No. of all possible actions

Policy iteration won't repeat policies.

14

Lemma: Let T be a γ -contraction mapping, $T: \mathbb{R}^N \rightarrow \mathbb{R}^N$

* * then $V_{n+1} = TV_n$

* then

$$\lim_{n \rightarrow \infty} V_n = V_{\gamma^*}$$

$$TV_{\gamma^*} = V_{\gamma^*}$$

Fixed point.
(Unique)

Show its cauchy,

$\Rightarrow \|V_{n+m} - V_n\| \leq \epsilon$ } For every $\epsilon > 0$, there exists n_0 such that $n, m \geq n_0$

$$\text{LHS: } \|V_{n+m} - V_n\| = \left\| \sum_{k=0}^{m-1} (V_{n+k+1} - V_{n+k}) \right\|$$

$$\leq \sum_{k=0}^{m-1} \|V_{n+k+1} - V_{n+k}\|$$

$$= \sum_{k=0}^{m-1} \|T V_{n+k} - T V_{n+k}\|$$

$$\leq \sum_{k=0}^{m-1} \gamma \|V_{n+k} - V_{n+k-1}\|$$

$$\leq \sum \gamma^{n+k} \|V_1 - V_0\|$$

$$= \sum_{k=0}^{m-1} \gamma^{n+k} \|V_1 - V_0\|$$

$$\text{G.P } \left\{ \begin{array}{l} = \gamma^n \frac{(1-\gamma^m)}{1-\gamma} \cdot \|V_1 - V_0\| \\ (\gamma < 1) \end{array} \right\}$$

$$\leq \epsilon$$

Since all
the terms
are
controllable
and they
don't diverge

$$\Rightarrow \lim_{n \rightarrow \infty} \gamma^n = 0, \text{ thus } \exists n_0(\epsilon)$$

$$\text{s.t. } \forall n \geq n_0, \|V_{n+m} - V_n\| \leq \epsilon$$

\Rightarrow Let V^* be the limit,

Thus, $\lim_{n \rightarrow \infty} T^n V_0 = V^*$

Now, we show that V^* is a fixed point which means $T V^* = V^*$

a) $\|T V^* - V^*\| \geq 0$

b) $\|T V^* - V^*\| \leq \|T V^* - V_n\| + \|V_n - V^*\|$
 $\leq \|T V^* - T V_{n-1}\| + \|V_{n-1} - V^*\|$
 $\leq \gamma \|V^* - V_{n-1}\| + \|V_n - V^*\|$

But $\|V_{n-1} - V^*\| \leq \epsilon \forall n \geq n_0(\epsilon)$

$\therefore \|T V^* - V^*\| \leq \gamma \epsilon + \epsilon = \epsilon(\gamma+1)$

$\|T V^* - V^*\| \leq \epsilon(\gamma+1)$

$\Rightarrow T V^* = V^* \quad (\because \text{It holds for any } \epsilon)$
 $\quad \quad \quad (\because \text{Norm is getting close to zero})$

Theorem: The policy iteration algo generates a sequence of policies with non-decreasing performance.

$$V^{\pi_{k+1}} \geq V^{\pi_k} \quad \text{(Sliding optimality)}$$

Proof: $V^{\pi_k} = F^{\pi_k} V^{\pi_k}$ } F is the bellman operator

$$\Rightarrow V^{\pi_k} = F^{\pi_k} V^{\pi_k} \\ \leq F V^{\pi_k}$$

V^{π_k} is the fixed point of F^{π_k}

$$= F^{\pi_{k+1}} V^{\pi_k}$$

$$\Rightarrow F^{\pi_{k+1}} V^{\pi_k} \leq (F^{\pi_{k+1}})^2 V^{\pi_k} \rightarrow \text{Monotonicity}$$

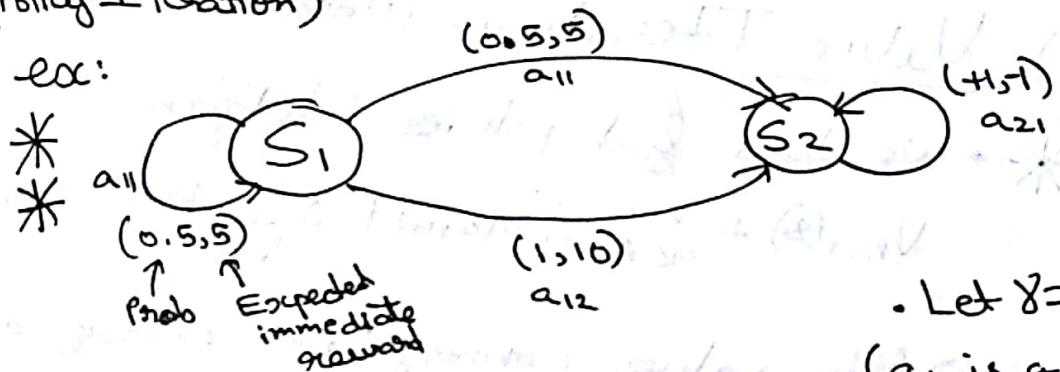
$$(F^{\pi_{k+1}})^{n-1} V^{\pi_k} \leq (F^{\pi_{k+1}})^n V^{\pi_k}$$

This will have a fixed point since it's a gamma contraction mapping

$$\Rightarrow V^{\pi_k} \leq F^{\pi_{k+1}} \cdot V^{\pi_k} \leq \dots \underbrace{(F^{\pi_{k+1}})^n V^{\pi_k}}_{V^{\pi_{k+1}}} \xrightarrow{\text{Fixed point of } F^{\pi_{k+1}}}$$

(Policy Iteration)

ex:



$$\cdot \text{Let } \gamma = 0.95$$

(a_{11} is a probabilistic action)

Step 1: $\pi^0(s_1) = a_{12}$ } Let this be our initial policy
 $\pi^0(s_2) = a_{21}$

• Policy evaluation: $V_{\pi^0}(s_1) = 10 + 0.95 [V_{\pi^0}(s_2)]$

$$V_{\pi^0}(s_2) = -1 + 0.95 [V_{\pi^0}(s_2)]$$

$$\therefore p(s_2|s_2, a_2) = 1$$

$$? V_{\pi^0}(s_1) = -9$$

$$V_{\pi^0}(s_2) = -20$$

• Policy improvement: $\pi^*(s_1) = \arg \max_{a \in A(s_1)} g_\gamma(s_1, a) + \gamma \sum_{s' \in A(s_1)}$

$$p(s'|s_1)$$

$$V_{\pi^0}(s_1)$$

$$= \max \left[\begin{array}{l} 5 + 0.95 \times \frac{1}{2} \times V_{\pi^0}(s_1) \\ 10 + 0.95 \times \frac{1}{2} \times V_{\pi^0}(s_2) \end{array} \right]$$

$$= \max \left[\begin{array}{l} -0.875 \\ -9 \end{array} \right]$$

$$\begin{aligned} \therefore \pi_1(s_1) &= a_{11} \\ \pi_1(s_2) &= a_{21} \end{aligned} \quad \left. \begin{array}{l} \text{After change of policy} \\ (s_2 \text{ has only 1 action}) \end{array} \right\}$$

Step 2:

- Policy evaluation: $V_{\pi_1}(s_1) = -8.571$
 $V_{\pi_1}(s_2) = -20$

- We keep repeating until $V_{\pi_{k+1}}(s)$ converges.
-

2) Value Iteration: (Details in slides)

- * * • We don't find policies in between

$$V_{n+1}(s) = \max_{a \in A(s)} \left[r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) \cdot V_n(s') \right]$$

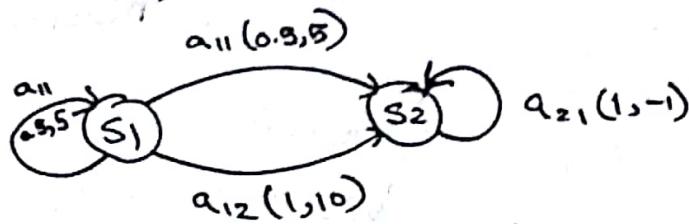
→ After values converge, we choose an output policy π^* ,

$$\pi^*(s) \in \arg \max_{a \in A(s)} \left[r(s, a) + \sum_{s' \in S} p(s'|s, a) \cdot V_{n+1}(s') \right]$$

- (15) • If $\|V_{n+1} - V_n\| < \epsilon \cdot \frac{1-\gamma}{2\gamma}$, then we stop.

(Value Iteration)

ex:



$$r_t(s, a, s') = \mathbb{E}[R_{t+1} | S_t = s, A_t = a, S_{t+1} = s']$$

$$\epsilon = 0.01, \gamma = 0.95$$

$$V_0(s_1) = V_0(s_2) = 0$$

$$V_{n+1}(s_1) = \max \left\{ \begin{array}{l} 5 + 0.5 \times 0.95 \times V_n(s_1) + 0.95 \times 0.95 \times V_n(s_2) \\ 10 + 0.95 \times V_n(s_2) \end{array} \right\}$$

$$V_{n+1}(s_2) = -1 + 0.95 \times V_n(s_2)$$

$$V_1(s_1) = \max(5, 10) = 10$$

$$V_2(s_2) = -1$$

$$\Rightarrow n=162, \|V^{163} - V^{162}\| < \frac{0.01 \times 0.05}{2 \times 0.95} = 0.00026$$

will $\{V_{162}(s_1) = \cancel{-10.5} \dots\}$, $\pi_c(s_1) = a_{11}$
converge $V_{162}(s_2) = \cancel{-10.5} \dots$, $\pi_c(s_2) = a_{21}$
to some values

Now consider optimality,

$$\Rightarrow V^*(s_1) = \max \begin{cases} 5 + 0.5 \times 8 \times V^*(s_2) + 0.15 \times 8 \times V_2(s_2) \\ 10 + 8 \times V_2(s_2) \end{cases}$$
$$V^*(s_2) = -1 + 8 \cdot V^*(s_2)$$

• Case 1: $\gamma = 0$
 $V^*(s_1) = 10$, $\pi^*(s_1) = a_{12}$
 $V^*(s_2) = -1$, $\pi^*(s_2) = a_{21}$

• Case 2: $\gamma = 0.5$

$$V^*(s_2) = -1 + 0.5 V_2(s_2) \Rightarrow V^*(s_2) = -2$$
$$V^*(s_1) = \max \begin{cases} 5 + 0.25 \times V^*(s_2) - 0.5 \\ 10 - 1 \end{cases}$$
$$= \max \{ 4.5 + 0.25 V^*(s_1) \}$$

$$V^*(s_1) \geq 4.5 + 0.25 V^*(s_1)$$

$$\Rightarrow 0.75 V^*(s_1) \geq 4.5$$

$$V^*(s_1) \geq 6 \quad 9 \text{ is better}$$

$$\therefore V^*(s_1) = 9$$

$$\pi^*(s_1) = a_{12}, \pi^*(s_2) = a_{21}$$

• Case 3: $\gamma = 0.95$

$$V^*(s_2) = -1 + 0.95 V_2(s_2)$$

$$V^*(s_2) = -20$$

$$V_{\pi^*}(s_1) = \max \left\{ \frac{5 + 0.475 V_{\pi^*}(s_1)}{10 - 20} \right.$$

$$= \max \left\{ \begin{array}{l} -4.5 + 0.475 \times V_{\pi^*}(s_1) \\ -10 \end{array} \right.$$

$$V_{\pi^*}(s_1) \geq -4.5 + 0.475 \times V_{\pi^*}(s_1)$$

$$V_{\pi^*}(s_1) \geq -\frac{4.5}{0.525}, \therefore V_{\pi^*}(s_1) \neq -10$$

$$\therefore V_{\pi^*}(s_1) = -\frac{4.5}{0.525}$$

$$\pi_{\pi^*}(s_1) = a_{11}, \pi_{\pi^*}(s_2) = a_{21}$$

* Correctness of Value Iteration Algo:

* Theorem: For the series V_n and policy π^* computed by VI, following will hold.

$$i) \lim_{n \rightarrow \infty} V_n = V_{\pi^*}$$

$$ii) \exists n_0 \text{ s.t. } \forall n \geq n_0, \|V_{n+1} - V_n\| \leq \epsilon - \frac{1-\gamma}{2\gamma}$$

* iii) The policy π^* is ϵ -optimal

$$\|V_{\pi^*} - V_{\pi^*}^{\pi^*}\| \leq \epsilon$$

$$iv) \text{ If } \|V_{n+1} - V_n\| < \epsilon - \frac{1-\gamma}{2\gamma} \text{ then }$$

$$\|V_{n+1} - V_{\pi^*}\| < \frac{\epsilon}{2}$$

(Assuming
h is
after
converged)

Proof:

* * iii) We have to show, $\|V_{\pi^*} - V_{\pi^*}^{\pi^*}\| \leq \epsilon$
LHS: $\|V_{\pi^*} - V_{\pi^*}^{\pi^*}\| = \|V_{\pi^*} - V_{n+1} + V_{n+1} - V_{\pi^*}^{\pi^*}\|$

D'le inequality \leftarrow $\leq \|V_{\pi^*} - V_{n+1}\| + \|V_{n+1} - V_{\pi^*}^{\pi^*}\|$

~~Part 1: $\|V_{\pi^*} - V_{n+1}\|$~~

~~$= \|V_{\pi^*} - FV_{n+1} + FV_{n+1} - V_{n+1}\|$~~

~~$\leq \|V_{\pi^*} - FV_{n+1}\| + \|FV_{n+1} - V_{n+1}\|$~~

$$\Rightarrow FV_{\gamma^*} = V_{\gamma^*}$$

$$\therefore \|V_{\gamma^*} - V_{n+1}\| \leq \|FV_{\gamma^*} - FV_{n+1}\|$$

$$+ \|FV_{n+1} - V_{n+1}\|$$

$$\leq \gamma \|V_{\gamma^*} - V_{n+1}\|$$

$$+ \gamma \|V_{n+1} - V_n\|$$

$$(1-\gamma) \|V_{\gamma^*} - V_{n+1}\| \leq \gamma \|V_{n+1} - V_n\|$$

$$\|V_{\gamma^*} - V_{n+1}\| \leq \frac{\gamma}{1-\gamma} \times \frac{\epsilon \times 1-\gamma}{2\gamma} = \frac{\epsilon}{2}$$

Part 2: $\|V_{\gamma^*}^{\pi^*} - V_{n+1}\| = \|V_{\gamma^*}^{\pi^*} - FV_{n+1}\| + FV_{n+1} - V_{n+1}\|$

$$\Rightarrow \|V_{\gamma^*}^{\pi^*} - V_{n+1}\| \leq \|V_{\gamma^*}^{\pi^*} - FV_{n+1}\| + \|FV_{n+1} - V_{n+1}\|$$

(Note: FV_{n+1}

$$= F^{\pi^*} V_{n+1})$$

Since At n, we
converge.

(Verify yourself)

$$\leq \|F_{V_{\gamma^*}}^{\pi^*} - F^{\pi^*} V_{n+1}\|$$

$$+ \|FV_{n+1} - FV_n\|$$

Similar to part 1,

$$\|V_{\gamma^*}^{\pi^*} - V_{n+1}\| \leq \frac{\epsilon}{2}$$

$$\therefore \|V_{\gamma^*}^{\pi^*} - V_{\gamma^*}^{\pi^*}\| \leq \epsilon$$

* Convergence of Value Iteration:

Claim: Monotonicity of convergence

if $FV^n \geq V^n$ then $\forall m \geq 0$

$$V_{n+m-1} \geq V_{n+m}$$

$$(F^{m+1}V^n \geq F^m V^n)$$

Theorem: Let V_n be the sequence of state values found by VI, then

$$i) \|V_n - V_{\gamma^*}\| \leq \frac{\gamma^n}{1-\gamma} \|V_1 - V_0\|$$

$$(ii) \|V_{\gamma^*} - V_{\gamma^*}^{\pi_e}\| \leq \frac{2\gamma^n}{1-\gamma} \|V_1 - V_0\|$$

$$\|V_{\gamma^*} - V_n\| \leq \frac{\gamma^n}{1-\gamma} \|V_n - V_0\| \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{We have shown this}$$

$$\leq \frac{\gamma^n}{1-\gamma} \|V_{n-1} - V_{n-2}\|$$

$$\leq \frac{\gamma^{2n}}{1-\gamma} \|V_{n-1} - V_{n-2}\|$$

$$\leq \frac{\gamma^n}{1-\gamma} \|V_1 - V_0\|$$

• Similarly, $\|V_{\gamma^*}^{\pi_e} - V_{\gamma^*}\| \leq \frac{\gamma^n}{1-\gamma} \|V_1 - V_0\|$

$$\begin{aligned} \|V_{\gamma^*}^{\pi_e} - V_{\gamma^*}\| &\leq \|V_{\gamma^*}^{\pi_e} - V_n + V_n - V_{\gamma^*}\| \\ &\leq \frac{2\gamma^n}{1-\gamma} \|V_1 - V_0\| \end{aligned}$$

(Convergence slower when γ is close to 1) \rightarrow Some improvements have been suggested

⑯ Monte Carlo Methods: (MDP₃)

- Monte Carlo methods learn from experience
- On-line experience
- Episodic tasks with reward at the end of each episode.
- * • Each occurrence of state s in an episode is called a visit.
- ↳ Naive approach:

$$\hat{V}_\pi(s) = \frac{1}{m} \sum_{i=1}^m r_i \quad \text{Return of } T_i \text{ (episode)}$$

For each episodes,
run π from s for
m times.

(Actions may be
different at each
episode)

* MC Approaches:

i) ↳ First visit monte carlo:

- Average rewards for first time s is visited in an episode $\rightarrow \hat{V}_\pi(s) = \frac{1}{m} \sum_{i=1}^m r_i(s, T_i)$

ii) ↳ Every visit monte carlo:

- Average rewards for every time s is visited in an episode.

(Assuming we get rewards between episodes as well)

$$\hat{V}_\pi(s) = \frac{1}{\sum_{i=1}^m N_i(s)} \sum_{i=1}^m \sum_{j=1}^{N_i(s)} r_i(s, T_{i,j})$$

Number of occasions
of s in an episode.

Note: example in slides — Blackjack — First visit MC

⇒ MC does not require full knowledge of the environment unlike DP.

- State value estimates for each state are independent unlike DP.



- Useful when underlying model is not available.
- For every state-action pair, performs mean reward calculation like earlier. (follow first)
 - ↳ we need to maintain exploration states visit MC so that we actually visit different state with different trials even in deterministic policies.

* MC Control: (Slides)

- Use estimation & improvement cycles as in policy iteration.
 - Policy improvement: $\pi(s) = \max_a V_\pi(s, a)$
 - * • We use exploration starts to visit all actions from the start state s to get a better approximation of $V_\pi(s, a)$ than.
 - Ideally we want ∞ episodes, but we can allow a tolerance in errors for finite episodes.

Without this we would only update $V_\pi(s, a_i)$ where $\pi(s) = a_i$, so we won't be able to update π properly.
- To fix this issue we follow method similar to ϵ -Greedy. (Best arm with ϵ , random with $1-\epsilon$ etc.)

END - Mid 2

- ⑦ * Issues with exploration starts is optimal not random
- * So we use ϵ -Greedy like approach.
- On we could use "Off Policy" where one policy is used as the behaviour policy and the policy being learnt is called the target policy. (ex: behaviour is stochastic and target is deterministic)
- (On Policy is evaluate and improve on the same policy)

Note: Example of first visit ~~the~~ monte carlo in slides.

- In "On Policy" methods we need non ~~deterministic~~ deterministic policies since we want to have exploration. (ex: ϵ -Greedy)

⇒ "On Policy" & "Off Policy" are two categories.

* 1) On Policy First Visit: (ϵ -Greedy like)

- After evaluation of $Q(s,a)$ get $A^* = \arg\max Q(s,a)$ and then set,

$$\pi(a|s) = 1 - \epsilon + \epsilon / |A(s)| \text{ if } a = A^*$$

$$\epsilon / |A(s)| \text{ if } a \neq A^*$$

- Evaluation is just average reward over all the runs.

$$q_{\pi}(s, \pi_1(s)) = q_{\pi}(s, A^*)(1 - \epsilon) + \sum_a \frac{\epsilon}{|A(s)|} q_{\pi}(s, a)$$

$$\geq (1 - \epsilon) \sum_a \left(\pi(a|s) - \frac{\epsilon}{|A(s)|} \right) q_{\pi}(s, a)$$

$$+ \frac{\epsilon}{|A(s)|} \sum_a q_{\pi}(s, a)$$

$$= \sum_a \pi(a|s) q_{\pi}(s, a) = v_{\pi}(s)$$

$\therefore q_{\pi}(s, \pi_1(s)) \geq v_{\pi}(s)$

∴ State values are \geq older ones.

State values

* 2) Off Policy Control:

- If policies are π & μ . Then if action has non-zero probability in π ~~has non-zero probability in π~~
 $\pi(a|s) > 0 \Rightarrow \mu(a|s) > 0$
→ Assumption of coverage

- Typically π would be a greedy policy.

$\Rightarrow \pi = \text{Target policy}$

$\lambda = \text{Behavioral policy}$

* Importance Sampling:

$$\mathbb{E}_{x \sim p}[\hat{s}] = \sum_x [p(x) \cdot \frac{x q(x)}{p(x)}] = \mathbb{E}_{x \sim q}[x]$$

$$(\hat{s} = \frac{x q(x)}{p(x)})$$

That is if $\hat{s} = \frac{q(s)}{p(s)}$

- We use importance sampling among π & λ .

- Probability of state-action trajectory under policy π is,

$$P_\pi(\{(s_t, a_t), \dots, (s_T, a_T)\})$$

$$= \prod_{k=1}^{T-1} \pi(a_k | s_k) \cdot p(s_{k+1} | s_k, a_k)$$

Since its markov

$$= \prod_{k=1}^{T-1} \pi(a_k | s_k) \cdot p(s_{k+1} | s_k, a_k)$$

Considering, $P_\pi^T = \frac{\prod_{k=1}^{T-1} \pi(a_k | s_k) \cdot p(s_{k+1} | s_k, a_k)}{\prod_{k=1}^{T-1} \lambda(a_k | s_k) \cdot p(s_{k+1} | s_k, a_k)}$

$$= \prod_{k=1}^{T-1} \pi(a_k | s_k)$$

$$\frac{\prod_{k=1}^{T-1} \pi(a_k | s_k)}{\prod_{k=1}^{T-1} \lambda(a_k | s_k)}$$

- $T(s) = \text{Set of all time steps where } s \text{ is visited. For first visit } T(s) \text{ is just a value of first time stamp.}$

- $T(t)$ is the first time of termination following time t .
 $G(t)$ is return after t till $T(t)$

$\Rightarrow \{G_t\}_{t \in T(s)}$ are returns of corresponding to state s . $\{\rho_t^{T(t)}\}_{t \in T(s)}$ are importance sampling ratios.

e.g.: $s_1, a_1, s_2, a_2, s_1, a_2, s_3, a_3, s_2, a_2, s_1$

First visit: $T(s_1) = 1$, $T(1) = 6$

Every visit: $T(s_1) = \{1, 3, 6\}$, $T(1) = 3$
 $T(3) = 6$

$$*\cdot \hat{v}(s) = \frac{\sum_{t \in T(s)} \rho_t^{T(t)} \cdot G_t}{|\Pi(s)|}$$

Importance sampling
↑ ↑ Returns of policy Π

This is ordinary importance sampling.

(8)

* \Rightarrow It is an unbiased estimate of $V_\pi(s)$

$$\mathbb{E}_\mu[\hat{v}_\pi(s)]$$

$$\text{Proof: } \mathbb{E}_\mu[\hat{v}_\pi(s)] = \mathbb{E}_\mu \left[\frac{\sum_{t \in T(s)} \rho_t^{T(t)} G_t}{|\Pi(s)|} \right]$$

$$= \frac{1}{|\Pi(s)|} \sum_{t \in T(s)} \mathbb{E}_\mu \left[\rho_t^{T(t)} G_t \right]$$

$$\cdot \mathbb{E}_\mu \left[\rho_t^{T(t)} G_t \right] = \mathbb{E}_\mu \left[\prod_{k=t}^{T(t)-1} \frac{\pi(A_k | S_k)}{M(A_k | S_k)} G_t \right]$$

$$= \mathbb{E}_\mu \left[\prod_{k=t}^{T(t)-1} \frac{\pi(A_k | S_k)}{M(A_k | S_k)} \cdot \sum_{R=0}^{T(t)-t-1} \gamma^k R \right]$$

$$= \sum_{R=0}^{T(t)-t-1} \gamma^k R \cdot \mathbb{E}_\mu \left[\prod_{k=t}^{T(t)-1} \frac{\pi(A_k | S_k)}{M(A_k | S_k)} \cdot R \right]$$

$$= \sum_{R=0}^{T(t)-t-1} \gamma^k R \left[\mathbb{E}_\mu \left(\frac{\pi(A_t | S_t)}{M(A_t | S_t)} \cdot R \right) \right] \prod_{k=t+1}^{T(t)-1} \mathbb{E}_\mu \left[\frac{\pi(A_k | S_k)}{M(A_k | S_k)} \right]$$

(The Π term
? is just 1)

$$= \sum_{R=0}^{T(t)-t-1} \gamma^k R \left[\mathbb{E}_\mu \left(\frac{\pi(A_t | S_t)}{M(A_t | S_t)} \cdot R \right) \right] \prod_{k=t+1}^{T(t)-1} \mathbb{E}_\mu \left[\frac{\pi(A_k | S_k)}{M(A_k | S_k)} \right]$$

$$= \sum_{t=0}^{T(s)-1} \gamma^t \mathbb{E}_{\pi}[R_{s+t+1}]$$

$$= \mathbb{E}_{\pi}[G_s | S_s = s] \quad \text{[Redacted]}$$

$$\therefore \mathbb{E}_{\mu}[\hat{v}_{\pi}(s)] = \mathbb{E}_{\pi}\left[\sum_{t=0}^{T(s)-1} \gamma^t R_{s+t+1} | S_s = s\right] = v_{\pi}(s)$$

Note:

$$Var(x) = \mathbb{E}[x^2] - [\mathbb{E}(x)]^2$$

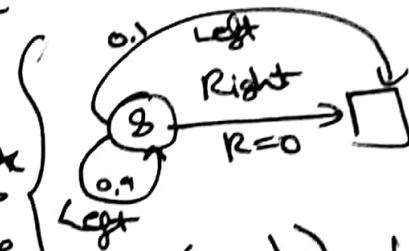
$\mathbb{E}[x^2]$ can make this unbounded since $\mathbb{E}(x)$ is bounded

ex:

$$\star \cdot \hat{v}_{\pi}(s) = \prod_{t=0}^{T-1} \frac{\pi(A_t | S_t)}{\mu(A_t | S_t)} \cdot G_0$$

$$\mathbb{E}_{\mu}[(\hat{v}_{\pi}(s))^2] = \mathbb{E}_{\pi}\left[\left(\prod_{t=0}^{T-1} \frac{\pi(A_t | S_t)}{\mu(A_t | S_t)} G_0\right)^2\right]$$

~~Consider trajectory:~~
 $s, left, 0, s, left, 0$
 $s, left, 0, s, left, +1$ Consider this example



$$\pi(left|s) = 1$$

$$\mu(left|s) = 1/2$$

$$\Rightarrow \mathbb{E}_{\mu}[(\hat{v}_{\pi}(s))^2] = \underbrace{\frac{1}{2} \times 0.1 \times \left(\frac{1}{0.5}\right)^2}_{\text{Length 1 trajectory}} + \underbrace{\frac{1}{2} \times 0.9 \times \left(\frac{1}{0.9}\right)^2}_{\text{Length 2 trajectory}} + \dots$$

$$= 0.1 \sum_{k=0}^{\infty} (0.9)^k \cdot 2^k \cdot 2 = \infty$$

(Variance is not bounded using ordinary importance sampling)

* Weighted Importance Sampling:

$$* \cdot E_{\pi}[\hat{v}_{\pi}(s)] = E_{\pi} \left[\frac{\sum_{t \in T(s)} p_t^{\pi} G_t}{\sum_{t \in T(s)} S_t^{\pi}} \right] = v_{\pi}(s)$$

\therefore This is not an unbiased estimator like earlier.

(Variance however converges to zero here)

Note: Off Policy every visit MC evaluation algo in slides.
 ** (Both with π given & unknown π) \rightarrow Initialized randomly & improved.

* Discounting Aware Importance Sampling:

- Why should G_{it} be weighted by $\prod_{k=0}^{T-1} \frac{\pi(A_k|S_k)}{\pi^*(A_k|S_k)}$, instead use $\frac{\pi(A_i|S_0)}{\pi^*(A_i|S_0)}$? \rightarrow Could have high variance

\Rightarrow Flat Partial Returns:

$$\bar{G}_{it} = \sum_{k=i}^h R_{i+k}, 0 \leq i < h \leq T$$

$$\begin{aligned} \bar{G}_{it} &= (1-\gamma)R_{i+1} + \gamma(1-\gamma)(R_{i+2} + R_{i+3}) \\ &\quad + \gamma^2(1-\gamma)(R_{i+4} + R_{i+5} + R_{i+6}) \\ &\quad + \dots \\ &\quad (1-\gamma)\gamma^{T-i-1}(R_{T-i} + \dots + R_T) \end{aligned}$$

$\gamma = \text{prob of termination at a given timestep}$

$$G_{it} = (1-\gamma) \sum_{h=i+1}^{T-1} \gamma^{h-i-1} \bar{G}_{it} + \gamma^{T-i-1} \bar{G}_{it}$$

• Ordinary importance sampling

$$\hat{v}(s) = \sum_{t=T(s)}^h \left\{ (1-\gamma) \sum_{h=t+1}^{T-1} \gamma^{h-t-1} \bar{G}_{it} + \gamma^{T-t-1} \bar{G}_{it} \right\}$$

$$\bar{S}_i^h = \prod_{k=i}^h \frac{\pi(A_k|S_k)}{\pi^*(A_k|S_k)} |T(s)|$$

Similarly we can define weighted sampling.

19

Partially observable MDPs:

- Not completely sure about which state it is in. set of observations observation probabilities

$$\langle S, A, T, R, \gamma, O \rangle$$

- We move from belief state b to b' on taking action a and getting an observation o .

$$\therefore b'(s_j) = P(s_j | o, a, b) = \frac{P(o | s_j, a) \sum_{s_i} P(s_i | s_j, a)}{\sum_{s_j} P(o | s_j, a) \sum_{s_i} P(s_i | s_j, a)}$$

↓
State estimator

* * * \Rightarrow Converting POMDPs to belief-MDPs:

$$\langle S, A, T, R, \gamma, O \rangle$$

↓

$$\langle B, A, P, R, \gamma \rangle \quad \{ \text{Belief MDP}$$

$$\cdot g(b, a) = \sum_{s \in S} b(s) \cdot R(s, a)$$

$$\cdot P(b' | b, a) = \sum_{o \in O} P(b' | b, a, o) \cdot P(o | a, b)$$

↳ Discretizing states:

- Belief state needs to be discretized, so that Belief MDP has finite number of states.
- This would however lead to an exponential number of states compared to POMDPs. → Practically infeasible generally

(we can use bellman equation to solve this)

* Policy Tree:

$$\star \cdot V_{\pi}(b) = \sum_{s_j} b(s_j) \cdot V_{\pi}(s_j) = \underbrace{b \cdot V_{\pi}}_{\alpha_{\pi} \text{ (Representation)}}$$

$\rightarrow V_t(b) = \max_{\pi} V_{\pi}(b)$



$$\cdot V_{\pi}(s) = R(s, a(\pi)) + \gamma \sum_{s' \in S} T(s'|s, a(\pi)) \cdot \sum_{a' \in A} T(s', a(\pi), o_i) \cdot V(s')$$

- This is similar to DP in MDPs, it is simpler to compute than $V_{\pi}(b)$ and hence can be used in practice.

$V_{\pi}(b)$ is linear in b , however taking $V_t(b)$ will be a combination of policies ~~for~~.
 (This shows that $V_t(b)$ is piecewise linear function & will be convex)
 \rightarrow (Not proved completely though)

Note: Witness algorithm

(20) TD Learning: (Temporal Difference)

- Apply updates at every epoch and not have to wait for an entire episode.

$$\cdot V(s_t) = V(s_t) + \alpha (r_{t+1} + \gamma V(s_{t+1}) - V(s_t))$$

A [↑] naive TD approach. $\rightarrow [TD(0)]$

Note: Sampling & Bootstrapping in DP, Monte Carlo & TD.

Note: TD(0) estimation of state values (in slides)
↳ Generates sequence of states in a row
& applies state value update considering
state s & next state s' .

- TD converges with probability 1 if we use decreasing alphas (with certain properties)
 - Batch updating in TD & MC: Average updates over a batch of episodes & update actual statevalues with these average

we can use ←
incremental
batch sizes.

$$\Rightarrow Q(S_t, A_t) = Q(S_t, A_t) + \alpha(R_{t+1} + \gamma \cdot Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$$

* ↑
Updating Q values
using TD -

* Sarsa: (On-Policy TD Control)

- Apply Q value update to greedily find best actions to update policy iteratively.
- ϵ -greedy to choose best action or random action.

$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha (R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$$

(21)

Note: GLIE (Greedy in the Limit with Infinite Exploration)

- ↳ G-Greedy falls in this
- ↳ All State action pairs visited ∞ times
- ↳ Converges to a solution.

- SARSA — State Action Reward State Action
- Also in slides — SARSA is GLIE sequence of policies
- SARSA converges if $\sum_{t=1}^{\infty} \alpha_t = \infty$ & $\sum_{t=1}^{\infty} \alpha_t^2 < \infty$

* Q - Learning: (Off policy TD control)

$$Q_t(S_t, A_t) = (1-\alpha) Q_{t-1}(S_t, A_t) + \alpha [R_{t+1} + \gamma V_{t-1}(S_{t+1})]$$

$$\text{where } V_{t-1}(S_{t+1}) = \max_a Q_{t-1}(S_{t+1}, a)$$

- This is off policy since updates use older policy and then the policy changes as we move on.

(Updates use $t-1$ th policy, whereas selection uses t th policy)

Since we do this it is different from SARSA where we use chosen action a' instead.

↑ Due to this it is off policy

- Q-Learning is more aggressive than SARSA since we use the max term.

\Rightarrow In examples we keep proceeding in episodes until we reach terminal states & then apply all the updates.

* Convergence of Q-Learning:

Define: Contraction operator H as follows,

$$H^Q(s, a) = \sum_{s' \in S} p(s'|s, a) [q_1(s', a) + \gamma \max_{a' \in A(s')} q_1(s', a')]$$

$$H: R^d \rightarrow R^d, d=|S| \times |A|$$

Lemma: H is a contraction operator in the ∞ norm

$$\|HQ_1 - HQ_2\|_\infty \leq \gamma \|Q_1 - Q_2\|_\infty$$

Proof:

$$\|HQ_1 - HQ_2\|_\infty = \max_{s, a} \left| \sum_{s' \in S} p(s'|s, a) [q_1(s', a) + \gamma \max_{a'} q_1(s', a')] - [q_2(s', a) + \gamma \max_{a'} q_2(s', a')] \right|$$

$$\text{Note: } \|v_1 - v_2\|_\infty = \max_{i \in [d]} |v_1^i - v_2^i|$$

$$\therefore \|HQ_1 - HQ_2\|_\infty = \max_{s, a} |HQ_1(s, a) - HQ_2(s, a)| \quad \text{①}$$

$$\text{From ①, } \|HQ_1 - HQ_2\|_\infty = \max_{s, a} \left| \gamma \sum_{s'} p(s'|s, a) \left(\max_{a'} Q_1(s', a') - \max_{a'} Q_2(s', a') \right) \right|$$

$$\leq \max_{s, a} \gamma \sum_{s'} \left| p(s'|s, a) \left(\max_{a'} Q_1(s', a') - \max_{a'} Q_2(s', a') \right) \right|$$

$$\leq \max_{s, a} \gamma \sum_{s'} p(s'|s, a) \max_{a'} |Q_1(s', a') - Q_2(s', a')|$$

$$\leq \max_{s, a} \gamma \sum_{s'} p(s'|s, a) \max_{s', a'} |Q_1(s', a') - Q_2(s', a')|$$

$$\leq \max_{s, a} \gamma \sum_{s'} p(s'|s, a) \|Q_1 - Q_2\|_\infty$$

$$= \underline{\gamma \|Q_1 - Q_2\|_\infty}$$

Theorem: The random process $\{\Delta_t\}$ taking values in \mathbb{R}^n and defined as,

$$\Delta_{t+1}(x) = (1 - \alpha_t(x))\Delta_t(x) + \alpha_t(x)$$

$$= \Delta_t F_t(x)$$

Prob
not
required
(Slide has
a slightly
different
version)

converges to 0 with prob 1 under the following assumptions,

Step
size
should
satisfy
these.

$$i) 0 \leq \alpha_t \leq 1, \sum_{t=1}^{\infty} \alpha_t = \infty, \sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

$$ii) \|\mathbb{E}[F_t(x) | \mathcal{F}_t]\|_W \leq \gamma \|\Delta_t\|_W$$

with $\gamma < 1$

$$iii) \text{Var}[F_t(x) | \mathcal{F}_t] \leq C [(\|\Delta_t\|_W)^2]$$

History till
time t

for $C > 0$

General
theorem
for stochastic
approximation

22

"Almost surely" — Law of large numbers

Converges with probability 1.

* Convergence proof of Q-Learning:

* Let $\{Q_t(s, a)\}_{t \geq 0}$ denote the sequence of state-action value functions generated by Q-learning algo.

$$Q_{t+1}(s, a) = Q_t(s, a) + \alpha_t(s, a) [g_t(s, a) + \gamma \max_{a'} Q_t(s', a') - Q_t(s, a)]$$

$$= Q_t(s, a) + \alpha_t(s, a) [g_t(s, a) + \gamma \mathbb{E}_{s' \sim P[s'|s, a]} (\max_{a'} Q_t(s', a')) - Q_t(s, a)]$$

$$= Q_t(s, a) + \alpha_t(s, a) [\gamma \max_{a'} Q_t(s', a') - Q_t(s, a)]$$

$$+ \alpha_t(s, a) [\gamma \max_{a'} Q_t(s', a') - \gamma \mathbb{E}_{s' \sim P[s'|s, a]} (\max_{a'} Q_t(s', a'))]$$

$$= Q_t(s, a) + \alpha_t(s, a) [\gamma \max_{a'} Q_t(s', a') - \gamma \mathbb{E}_{s' \sim P[s'|s, a]} (\max_{a'} Q_t(s', a'))]$$

$$(Q_t \in \mathbb{R}^d)$$

* $Q_t(s, a)$ defines a component corresponding to $s_t = s$, $A_t = a$.

$$W_t(s) = \gamma \max_{a'} Q_t(s', a') - \gamma \mathbb{E}_{s' \sim P_t(s'|s, a)} \max_{a'}$$

Now,

$$\mathbb{E}[W_t | f_t] = 0, \text{ so this assumption of general theorem of stochastic approximation is satisfied.}$$

$$\Rightarrow Q_{t+1}(s, a) = Q_t(s, a) + \alpha_t(s, a) [HQ_t(s, a) - W_t(s)] \quad \rightarrow ①$$

$$(\because HQ(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P_t(s'|s, a)} \max_{a'} Q_t(s', a'))$$

We can see that ① represents the formulation as in general theorem for stochastic approx.

$$\Rightarrow |W_t(s)| \leq \gamma \max_{a'} |Q_t(s', a')| + \gamma |\mathbb{E}_{s' \sim P_t(s'|s, a)} \max_{a'} Q_t(s', a')|$$

$$\leq \gamma \max_{a', s'} |Q_t(s', a')| + \gamma \cdot \max_{a', s'} |Q_t(s', a')|$$

$$= 2\gamma \max_{a', s'} |Q_t(s', a)| = 2\gamma \|Q_t\|_\infty$$

∴ All conditions are met.

∴ It converges.

* Expected SARSA:

$$\text{Sarsa: } Q_t(s, a) = Q_t(s, a) + \alpha_t(s, a) [R_{t+1}(s, a) + \gamma Q_{t+1}(s', a') - Q_t(s, a)]$$

Expected Sarsa: ~~Sarsa~~

$$Q_t(s, a) = (1-\alpha)Q_{t-1}(s, a)$$

$$+ \alpha \left[R_{t+1} + \gamma \sum_{a' \in A} \pi(a'|s) \cdot Q_{t+1}(s', a') \right]$$

We use probabilities of all actions, unlike Sarsa where mostly the best action only contributed.

mostly the best action only contributed.

- We can prove convergence of expected SARSA in a manner similar to convergence of Q-Learning.
 - Performs better than SARSA but is computationally more expensive.
-

* Maximization Bias:

- In Q-Learning & SARSA, max over estimated values is used implicitly as an estimate of the maximum value which leads to a positive bias.

(Since observed values \rightarrow max over them is not the mean, but higher)

↳ In Q-Learning, early on we choose a lot of wrong actions but later on it's better. ~~but~~ Those wrong actions are due to maximization bias.

* Double Q-Learning: (Algo in slides)

* Use 2 action-value functions Q_1 & Q_2

- Do Q learning on both independently, randomly updating Q_1 or Q_2 at each step.
 - On updating Q_1 , use Q_2 for the value of the next state. Max action over Q_1 only however.
 - Action selections are ϵ greedy w.r.t sum of Q_1 & Q_2 .
-

* Multistep Bootstrapping:

- Updates are not performed at every timestep, but uses a few timesteps together.
- In between TD(0) & Monte-Carlo.

② Multistep TD learning:

* Infinite step TD is monte carlo
 ↳ 1 update per episode

⇒ n step return:

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} V_{t+n}^{(n)}$$

? ⇒ $V_{t+n}^{(n)} = V_{t+n-1}^{(n)} + \alpha [G_t^{(n)} - V_{t+n-1}^{(n)}]$
 (Updates start after $n-1$ steps)

→ To make up for those $n-1$ updates we perform $n-1$ additional updates at the end of the episode.

- In practice from $n=1$ ($\text{TD}(0)$), as we increase n , MSE drops but after a while it starts to rise again.
- Step size α also has an optimum inbetween where it increases as we move away from it.

* ⇒ Error reduction property: (Used to show convergence of multi-step TD)

$$\max_s |E_\pi[G_t^{(n)} | S_t = s] - V_\pi(s)| \leq \gamma^n \max_s |V_{t+n-1}(s) - V_\pi(s)|$$

Proof:

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} V_{t+n}^{(n)}$$

$$= G_t + \sum_{k=n}^{\infty} \gamma^k R_{t+k-1} + \gamma^{n-1} V_{t+n-1}^{(n)}$$

$$\Rightarrow E_\pi[G_t^{(n)} | S_t = s] - E_\pi[G_t | S_t = s]$$

$$= E_\pi[\gamma^{n-1} V_{t+n-1}^{(n)}] - \sum_{k=n}^{\infty} \gamma^k R_{t+k-1}$$

$$= \mathbb{E}_\pi \left[\gamma^n V_{t+n}(\mathbf{s}_{t+n}) - \gamma^n G_{t+n} \right]$$

$$\therefore \gamma^n \sum_{k=n}^{\infty} \gamma^{k-n} R_{t+k+1} = \sum_{l=0}^{\infty} R_{t+l+1} = \gamma^n G_{t+n}$$

$$= \gamma^n G_{t+n}$$

$$= \gamma^n \mathbb{E}_\pi \left[V_{t+n-1}(\mathbf{s}_{t+n}) - G_{t+n} \right]$$

$$\mathbb{E}_\pi \left[G_{t+n} | \mathbf{s}_t = s \right] - V_\pi(s) \xrightarrow{\mathbb{E}_\pi[G_{t+n}] = V_\pi(s)}$$

$$= \mathbb{E}_\pi \left[\gamma^n V_{t+n-1}(\mathbf{s}_{t+n}) - \gamma^n G_{t+n} \right]$$

$$= \mathbb{E}_{A_{t+n} | \mathbf{s}_{t+n}} \left[\mathbb{E}_\pi \left[(\gamma^n V_{t+n-1}(\mathbf{s}_{t+n}) - \gamma^n G_{t+n}) | \mathbf{s}_{t+n} \right] \right]$$

$$= \mathbb{E}_{\cancel{A_{t+n}} | \mathbf{s}'_t} \left[\gamma^n V_{t+n-1}(s') - \gamma^n V_\pi(s') \right]$$

$$\therefore \left| \mathbb{E}_\pi \left[G_{t+n} | \mathbf{s}_t = s \right] - V_\pi(s) \right|$$

$$= \left| \mathbb{E}_{s'} \left[\gamma^n V_{t+n-1}(s') - \gamma^n V_\pi(s') \right] \right|$$

$$\leq \mathbb{E}_{s'} \left| \gamma^n V_{t+n-1}(s') - \gamma^n V_\pi(s') \right|$$

$$\leq \max_{s'} \left| \gamma^n V_{t+n-1}(s') - \gamma^n V_\pi(s') \right|$$

Since this holds for all s ,

$$\max_s \left| \mathbb{E}_\pi \left[G_{t+n} | \mathbf{s}_t = s \right] - V_\pi(s) \right|$$

$$\leq \max_s \left| \gamma^n V_{t+n-1}(s) - V_\pi(s) \right|$$

Note: n-step SARSA & n-step expected SARSA
are done in a similar manner (skipping)

⇒ n-Step Off Policy by Importance Sampling



$$g_{t+n} = \prod_{k=t}^{\min(t+n-1, T-1)} \frac{\pi(A_k | S_k)}{\pi'(A_k | S_k)}$$



Off policy
n-step TD

$$V_{t+n}(s_t) = V_{t+n-1}(s_t) + \alpha g_{t+n}(G_t^n) - V_{t+n-1}(s_t)$$

Off policy
n-step SARSA

$$Q_{t+n}(s_t, a_t) = Q_{t+n-1}(s_t, a_t) + \alpha g_{t+n}(G_t^n) - Q_{t+n-1}(s_t, a_t)$$

(G_t^n would vary
for expected SARSA)

Note: n-Step tree backup algorithm