Learning to Optimize Exploration and Generalization

Benjamin Van Roy

work done with Dan Russo













exploration versus exploitation



$$\max_{x \in \mathbb{X}} f_{\theta}(x) \qquad \quad \theta \in \Theta$$

$$\max_{x \in \mathbb{X}} f_{\theta}(x) \qquad \quad \theta \in \Theta$$

• Expected reward $f_{\theta}(X_t) = E[R_t | X_t, \theta]$

$$\max_{x \in \mathbb{X}} f_{\theta}(x) \qquad \quad \theta \in \Theta$$

- Expected reward $f_{\theta}(X_t) = E[R_t | X_t, \theta]$
- Represent knowledge about model via

$$\max_{x \in \mathbb{X}} f_{\theta}(x) \qquad \quad \theta \in \Theta$$

- Expected reward $f_{\theta}(X_t) = E[R_t | X_t, \theta]$
- Represent knowledge about model via
 - Set membership $\theta \in \Theta_t \subseteq \Theta$

$$\max_{x \in \mathbb{X}} f_{\theta}(x) \qquad \quad \theta \in \Theta$$

- Expected reward $f_{\theta}(X_t) = E[R_t | X_t, \theta]$
- Represent knowledge about model via
 - Set membership $\theta \in \Theta_t \subseteq \Theta$
 - Probability distribution $\theta \sim p_t(\cdot)$







• Action/arm $\mathbb{X} = \{1, \dots, n\}$



- Action/arm $\mathbb{X} = \{1, \dots, n\}$
- Mean rewards with independent priors

$$f_{\theta}(x) = \theta_x$$
 $\theta \sim p_0(\theta) = \prod_{x=1}^N p_0^x(\theta_x)$



- Action/arm $\mathbb{X} = \{1, \dots, n\}$
- Mean rewards with independent priors

$$f_{\theta}(x) = \theta_x$$
 $\theta \sim p_0(\theta) = \prod_{x=1}^N p_0^x(\theta_x)$

• Feedback/reward $R_t = Y_t = f_{\theta}(X_t) + W_t$



- Action/arm $\mathbb{X} = \{1, \dots, n\}$
- Mean rewards with independent priors

$$f_{\theta}(x) = \theta_x$$
 $\theta \sim p_0(\theta) = \prod_{x=1}^N p_0^x(\theta_x)$

- Feedback/reward $R_t = Y_t = f_{\theta}(X_t) + W_t$
- Discounted objective addressed by Gittin's Index Theorem

• Linear program

$$f_{\theta}(x) = \theta^{\top} x$$
$$\mathbb{X} = \{x : Ax \le b\}$$

• Linear program

$$f_{\theta}(x) = \theta^{\top} x$$
$$\mathbb{X} = \{x : Ax \le b\}$$

Gaussian Prior

$$\theta \sim N(\mu, \Sigma)$$

• Linear program $f_{\theta}(x) = \theta^{\top} x$

$$\mathbb{X} = \{x : Ax \le b\}$$

- Gaussian Prior $\theta \sim N(\mu, \Sigma)$
- Noisy feedback / reward $R_t = Y_t = \theta^\top X_t + W_t$ $W_t \sim N(0, \sigma^2)$

• Linear program $f_{\theta}(x) = \theta^{\top} x$

$$\mathbb{X} = \{x : Ax \le b\}$$

- Gaussian Prior $\theta \sim N(\mu, \Sigma)$
- Noisy feedback / reward $R_t = Y_t = \theta^\top X_t + W_t$ $W_t \sim N(0, \sigma^2)$

natural objectives are intractable

• Comparisons via

- Comparisons via
 - Simulations

- Comparisons via
 - Simulations
 - Theoretical objectives such as expected regret

$$E\left[\operatorname{Regret}(T)\right] = \sum_{t=1}^{T} E\left[\max_{x} f_{\theta}(x) - f_{\theta}(X_{t})\right]$$

- Comparisons via
 - Simulations
 - Theoretical objectives such as expected regret

$$E\left[\operatorname{Regret}(T)\right] = \sum_{t=1}^{T} E\left[\max_{x} f_{\theta}(x) - f_{\theta}(X_{t})\right]$$
for optimal action

- Comparisons via
 - Simulations
 - Theoretical objectives such as expected regret

$$E\left[\operatorname{Regret}(T)\right] = \sum_{t=1}^{T} E\left[\max_{x} f_{\theta}(x) - f_{\theta}(X_{t})\right]$$
for optimal for selected action for selected action

- Comparisons via
 - Simulations
 - Theoretical objectives such as expected regret

$$E\left[\operatorname{Regret}(T)\right] = \sum_{t=1}^{T} E\left[\max_{x} f_{\theta}(x) - f_{\theta}(X_{t})\right]$$
expectation
over models
for optimal
action
for selected
action

- Comparisons via
 - Simulations
 - Theoretical objectives such as expected regret

$$E\left[\operatorname{Regret}(T)\right] = \sum_{t=1}^{T} E\left[\max_{x} f_{\theta}(x) - f_{\theta}(X_{t})\right]$$
expectation over models
for optimal action

minimizing expected regret maximizes expected reward

- Comparisons via
 - Simulations
 - Theoretical objectives such as expected regret

$$E\left[\operatorname{Regret}(T)\right] = \sum_{t=1}^{T} E\left[\max_{x} f_{\theta}(x) - f_{\theta}(X_{t})\right]$$
expectation
over models
for optimal
action
for selected
action

minimizing expected regret maximizes expected reward

• Emphasis has been on "large" *T*

- Comparisons via
 - Simulations
 - Theoretical objectives such as expected regret

$$E\left[\operatorname{Regret}(T)\right] = \sum_{t=1}^{T} E\left[\max_{x} f_{\theta}(x) - f_{\theta}(X_{t})\right]$$
expectation over models
for optimal action

minimizing expected regret maximizes expected reward

- Emphasis has been on "large" T
- Popular approaches to heuristic design

- Comparisons via
 - Simulations
 - Theoretical objectives such as expected regret

$$E\left[\operatorname{Regret}(T)\right] = \sum_{t=1}^{T} E\left[\max_{x} f_{\theta}(x) - f_{\theta}(X_{t})\right]$$
expectation
over models
for optimal
action
for selected
action

minimizing expected regret maximizes expected reward

- Emphasis has been on "large" T
- Popular approaches to heuristic design
 - Upper-confidence-bounds [Lai-Robbins, 1985; Dani-Hayes-Kakade, 2008; Rusmevichientong-Tsitsiklis, 2010; etc.]

- Comparisons via
 - Simulations
 - Theoretical objectives such as expected regret



minimizing expected regret maximizes expected reward

- Emphasis has been on "large" *T*
- Popular approaches to heuristic design
 - Upper-confidence-bounds [Lai-Robbins, 1985; Dani-Hayes-Kakade, 2008; Rusmevichientong-Tsitsiklis, 2010; etc.]
 - Thompson sampling [Thompson, 1933]

Upper-Confidence-Bound Algorithms (UCB)

Upper-Confidence-Bound Algorithms (UCB)

• Confidence set Θ_t
- Confidence set Θ_t
 - Set of "statistically plausible" models

- Confidence set Θ_t
 - Set of "statistically plausible" models
 - Updated based on observations

- Confidence set Θ_t
 - Set of "statistically plausible" models
 - Updated based on observations
- Upper confidence bounds $U_t(x) = \max_{\theta \in \Theta_t} f_{\theta}(x)$

- Confidence set Θ_t
 - Set of "statistically plausible" models
 - Updated based on observations
- Upper confidence bounds $U_t(x) = \max_{\theta \in \Theta_t} f_{\theta}(x)$
- Optimistic optimization $X_t \in \arg \max_{x \in \mathbb{X}} U_t(x)$

- Confidence set Θ_t
 - Set of "statistically plausible" models
 - Updated based on observations
- Upper confidence bounds $U_t(x) = \max_{\theta \in \Theta_t} f_{\theta}(x)$
- Optimistic optimization $X_t \in \arg \max_{x \in \mathbb{X}} U_t(x)$
- Bayes-UCB

- Confidence set Θ_t
 - Set of "statistically plausible" models
 - Updated based on observations
- Upper confidence bounds $U_t(x) = \max_{\theta \in \Theta_t} f_{\theta}(x)$
- Optimistic optimization $X_t \in \arg \max_{x \in \mathbb{X}} U_t(x)$
- Bayes-UCB
 - Maintain probability distribution $p_t(d\theta) = \mathbb{P}\left[\theta \in d\theta | \mathbb{F}_{t-1}\right]$

- Confidence set Θ_t
 - Set of "statistically plausible" models
 - Updated based on observations
- Upper confidence bounds $U_t(x) = \max_{\theta \in \Theta_t} f_{\theta}(x)$
- Optimistic optimization $X_t \in \arg \max_{x \in \mathbb{X}} U_t(x)$
- Bayes-UCB
 - Maintain probability distribution $p_t(d\theta) = \mathbb{P}\left[\theta \in d\theta | \mathbb{F}_{t-1}\right]$
 - Select level set as confidence set

• Maintain probability distribution $p_t(d\theta) = \mathbb{P}\left[\theta \in d\theta | \mathbb{F}_t\right]$

- Maintain probability distribution $p_t(d\theta) = \mathbb{P}\left[\theta \in d\theta | \mathbb{F}_t\right]$
- Sample model $\theta_t \sim p_t$

- Maintain probability distribution $p_t(d\theta) = \mathbb{P}\left[\theta \in d\theta | \mathbb{F}_t\right]$
- Sample model $\theta_t \sim p_t$
- Optimize sample $X_t \in \arg \max_{x \in \mathbb{X}} f_{\theta_t}(x)$

- Maintain probability distribution $p_t(d\theta) = \mathbb{P}\left[\theta \in d\theta | \mathbb{F}_t\right]$
- Sample model $\theta_t \sim p_t$
- Optimize sample $X_t \in \arg \max_{x \in \mathbb{X}} f_{\theta_t}(x)$

sample each action with the probability that it is optimal

- UCB
 - Finite indep. X

[Auer et al, 2002]

- UCB
 - Finite indep. X
 - Linear

[Auer et al, 2002] [Dani-Hayes-Kakade, 2008; Rusmevichientong-Tsitsiklis, 2010; Abbasi-Yadkori et al, 2011]

- UCB
 - Finite indep. X
 - Linear
 - Generalized linear

[Auer et al, 2002] [Dani-Hayes-Kakade, 2008; Rusmevichientong-Tsitsiklis, 2010; Abbasi-Yadkori et al, 2011]

[Filippi et al, 2010]

• UCB

TS

- Finite indep. X
- Linear
- Generalized linear

[Auer et al, 2002] [Dani-Hayes-Kakade, 2008; Rusmevichientong-Tsitsiklis, 2010; Abbasi-Yadkori et al, 2011]

[Filippi et al, 2010]

- UCB
 - Finite indep. X
 - Linear
 - Generalized linear
- [Auer et al, 2002] [Dani-Hayes-Kakade, 2008; Rusmevichientong-Tsitsiklis, 2010; Abbasi-Yadkori et al, 2011]
- [Filippi et al, 2010]

- TS
 - Finite X, Bernoulli [Agrawal-Goyal, 2012]

- UCB
 - Finite indep. X
 - Linear
 - Generalized linear
- [Auer et al, 2002] [Dani-Hayes-Kakade, 2008; Rusmevichientong-Tsitsiklis, 2010; Abbasi-Yadkori et al, 2011]
- [Filippi et al, 2010]

- TS
 - Finite X, Bernoulli
 - Linear

[Agrawal-Goyal, 2012]

[Agrawal-Goyal, 2012]

- UCB
 - Finite indep. X
 - Linear
 - Generalized linear

[Auer et al, 2002] [Dani-Hayes-Kakade, 2008; Rusmevichientong-Tsitsiklis, 2010; Abbasi-Yadkori et al, 2011]

[Filippi et al, 2010]

- \rightarrow TS
 - Finite X, Bernoulli [Agrawal-Goyal, 2012]
 - Linear

[Agrawal-Goyal, 2012]

UCB Regret Bounds — TS E[Regret] Bounds

[Russo-Van Roy, 2013]

- UCB
 - Finite indep. X
 - Linear
 - Generalized linear

[Auer et al, 2002] [Dani-Hayes-Kakade, 2008; Rusmevichientong-Tsitsiklis, 2010; Abbasi-Yadkori et al, 2011]

[Filippi et al, 2010]

- TS
 - Finite X, Bernoulli [Agrawal-Goyal, 2012]
 - Linear

[Agrawal-Goyal, 2012]

UCB Regret Bounds — TS E[Regret] Bounds

[Russo-Van Roy, 2013]

• The role of confidence sets

- UCB
 - Finite indep. X
 - Linear
 - Generalized linear

[Auer et al, 2002] [Dani-Hayes-Kakade, 2008; Rusmevichientong-Tsitsiklis, 2010; Abbasi-Yadkori et al, 2011]

[Filippi et al, 2010]

- TS
 - Finite X, Bernoulli
 - Linear

[Agrawal-Goyal, 2012]

[Agrawal-Goyal, 2012]

UCB Regret Bounds — TS E[Regret] Bounds

[Russo-Van Roy, 2013]

- The role of confidence sets
 - UCB: algorithm design and analysis
 - TS: analysis only

Learning to Optimize

Linear Bandit Simulations

$$f_{\theta}(x) = \theta^{\top} x \qquad \mathbb{X} = \{z_1, \dots, z_n\}$$



Learning to Optimize



Learning to Optimize



• TS is often tractable when Bayes-UCB is not

TS is often tractable when Bayes-UCB is not
Consider LP

$$f_{\theta}(x) = \theta^{\top} x \qquad \qquad \theta \sim N(\mu, \Sigma)$$
$$\mathbb{X} = \{x : Ax \le b\} \qquad \qquad R_t = Y_t = \theta^{\top} X_t + W_t$$
$$W_t \sim N(0, \sigma^2)$$

TS is often tractable when Bayes-UCB is not
Consider LP

$$f_{\theta}(x) = \theta^{\top} x \qquad \qquad \theta \sim N(\mu, \Sigma)$$
$$\mathbb{X} = \{x : Ax \le b\} \qquad \qquad R_t = Y_t = \theta^{\top} X_t + W_t$$
$$W_t \sim N(0, \sigma^2)$$

• TS is computationally efficient

TS is often tractable when Bayes-UCB is not
Consider LP

$$f_{\theta}(x) = \theta^{\top} x \qquad \qquad \theta \sim N(\mu, \Sigma)$$
$$\mathbb{X} = \{x : Ax \le b\} \qquad \qquad R_t = Y_t = \theta^{\top} X_t + W_t$$
$$W_t \sim N(0, \sigma^2)$$

- TS is computationally efficient
- Bayes-UCB is computationally intractable

TS is often tractable when Bayes-UCB is not
Consider LP

$$f_{\theta}(x) = \theta^{\top} x \qquad \qquad \theta \sim N(\mu, \Sigma)$$
$$\mathbb{X} = \{x : Ax \le b\} \qquad \qquad R_t = Y_t = \theta^{\top} X_t + W_t$$
$$W_t \sim N(0, \sigma^2)$$

- TS is computationally efficient
- Bayes-UCB is computationally intractable
- Computationally tractable version of UCB
 - Regret scaled by a factor of d [Dani-Hayes-Kakade, 2008]

Learning to Optimize

• TS outperforms Bayes-UCB designed for analysis

- TS outperforms Bayes-UCB designed for analysis
- TS slightly underperforms well-tuned Bayes-UCB

- TS outperforms Bayes-UCB designed for analysis
- TS slightly underperforms well-tuned Bayes-UCB
- TS often tractable when Bayes-UCB not
Summary of TS versus UCB

- TS outperforms Bayes-UCB designed for analysis
- TS slightly underperforms well-tuned Bayes-UCB
- TS often tractable when Bayes-UCB not
- TS outperforms non-Bayes-UCB designed for tractability

$$E\left[Regret(T)\right] \leq \tilde{O}\left(\sqrt{d_E(T)\log\left(N(T)\right)T}\right)$$
 [Russo-Van Roy, 2013]
$$T^{-2}\text{-scale eluder dimension}$$
 of function class of function class

• Bound via general notion of function class complexity

$$E\left[Regret(T)\right] \leq \tilde{O}\left(\sqrt{d_E(T)\log\left(N(T)\right)T}\right)$$
 [Russo-Van Roy, 2013]
$$T^{-2}\text{-scale eluder dimension}$$
 of function class of function class

• CN is representative of supervised learning concepts

$$E\left[Regret(T)\right] \leq \tilde{O}\left(\sqrt{d_E(T)\log\left(N(T)\right)T}\right) \qquad \text{[Russo-Van Roy, 2013]}$$

$$T^{-2}\text{-scale eluder dimension} \qquad T^{-2}\text{-covering number} \qquad \text{of function class} \qquad \text{of function class}$$

- CN is representative of supervised learning concepts
- ED is new and necessary

$$E\left[Regret(T)\right] \leq \tilde{O}\left(\sqrt{d_E(T)\log\left(N(T)\right)T}\right) \qquad \text{[Russo-Van Roy, 2013]}$$

$$T^{-2}\text{-scale eluder dimension} \qquad T^{-2}\text{-covering number} \qquad \text{of function class} \qquad \text{of function class}$$

- CN is representative of supervised learning concepts
- ED is new and necessary
- Specializes to various model classes

$$E\left[Regret(T)\right] \leq \tilde{O}\left(\sqrt{d_E(T)\log\left(N(T)\right)T}\right) \qquad \text{[Russo-Van Roy, 2013]}$$

$$T^{-2}\text{-scale eluder dimension} \qquad T^{-2}\text{-covering number} \qquad \text{of function class} \qquad \text{of function class}$$

- CN is representative of supervised learning concepts
- ED is new and necessary
- Specializes to various model classes
 - Linear bandits: recovers best previous bounds

$$E\left[Regret(T)\right] \leq \tilde{O}\left(\sqrt{d_E(T)\log\left(N(T)\right)T}\right) \qquad \text{[Russo-Van Roy, 2013]}$$

$$T^{-2}\text{-scale eluder dimension} \qquad T^{-2}\text{-covering number} \qquad \text{of function class} \qquad \text{of function class}$$

- CN is representative of supervised learning concepts
- ED is new and necessary
- Specializes to various model classes
 - Linear bandits: recovers best previous bounds
 - Generalized linear bandits: slight improvement

$$f_{\theta}(x) = \theta^{\top} x$$

$$f_{\theta}(x) = \theta^{\top} x \qquad \mathbb{X} = \left\{ x \in \left\{ 0, \frac{1}{m} \right\}^d : \|x\|_0 = m \right\}$$

$$f_{\theta}(x) = \theta^{\top} x \qquad \mathbb{X} = \left\{ x \in \left\{ 0, \frac{1}{m} \right\}^d : \|x\|_0 = m \right\}$$

$$R_t = Y_t = f_\theta(X_t)$$

$$f_{\theta}(x) = \theta^{\top} x \qquad \mathbb{X} = \left\{ x \in \left\{ 0, \frac{1}{m} \right\}^{d} : \|x\|_{0} = m \right\}$$
$$R_{t} = Y_{t} = f_{\theta}(X_{t})$$
$$\theta \sim \operatorname{unif}\left(\left\{ \theta \in \{0, 1\}^{d} : \|\theta\|_{0} = 1 \right\} \right)$$

• A 1-sparse case

$$f_{\theta}(x) = \theta^{\top} x \qquad \mathbb{X} = \left\{ x \in \left\{ 0, \frac{1}{m} \right\}^d : \|x\|_0 = m \right\}$$
$$R_t = Y_t = f_{\theta}(X_t)$$

$$\theta \sim \operatorname{unif}\left(\left\{\theta \in \{0,1\}^d : \|\theta\|_0 = 1\right\}\right)$$

• UCB/TS require $\Omega(d)$ samples to identify

$$f_{\theta}(x) = \theta^{\top} x \qquad \mathbb{X} = \left\{ x \in \left\{ 0, \frac{1}{m} \right\}^d : \|x\|_0 = m \right\}$$

$$R_t = Y_t = f_\theta(X_t)$$

$$\theta \sim \operatorname{unif}\left(\left\{\theta \in \{0,1\}^d : \|\theta\|_0 = 1\right\}\right)$$

- UCB/TS require $\Omega(d)$ samples to identify
 - UCB/TS rule out one action per period

$$f_{\theta}(x) = \theta^{\top} x \qquad \mathbb{X} = \left\{ x \in \left\{ 0, \frac{1}{m} \right\}^d : \|x\|_0 = m \right\}$$

$$R_t = Y_t = f_\theta(X_t)$$

$$\theta \sim \operatorname{unif}\left(\left\{\theta \in \{0,1\}^d : \|\theta\|_0 = 1\right\}\right)$$

- UCB/TS require $\Omega(d)$ samples to identify
 - UCB/TS rule out one action per period
- Easy to design algorithms for which $\log_2(d)$ suffice

• A simple context

- A simple context
 - *N* customer types

- A simple context
 - *N* customer types
 - Many products, each geared for a particular type

- A simple context
 - *N* customer types
 - Many products, each geared for a particular type
 - Action: recommend assortment of size M

- A simple context
 - *N* customer types
 - Many products, each geared for a particular type
 - Action: recommend assortment of size M
 - Customer purchases at most one product per period

- A simple context
 - *N* customer types
 - Many products, each geared for a particular type
 - Action: recommend assortment of size M
 - Customer purchases at most one product per period
 - Learn about customer through repeated interactions

- A simple context
 - *N* customer types
 - Many products, each geared for a particular type
 - Action: recommend assortment of size *M*
 - Customer purchases at most one product per period
 - Learn about customer through repeated interactions
- UCB/TS focus on a single customer type

- A simple context
 - *N* customer types
 - Many products, each geared for a particular type
 - Action: recommend assortment of size M
 - Customer purchases at most one product per period
 - Learn about customer through repeated interactions
- UCB/TS focus on a single customer type
- Diversifying can reduce regret by a factor of *M*

• Information ratio (IR)

$$\Psi_t = \frac{\left(E\left[f_{\theta}(X^*) - f_{\theta}(X_t) | \mathbb{F}_{t-1}\right]\right)^2}{I_t(X^*, Y_t)} = \frac{(\text{expected regret})^2}{\text{mutual information}}$$

• Information ratio (IR)

$$\Psi_t = \frac{\left(E\left[f_{\theta}(X^*) - f_{\theta}(X_t) | \mathbb{F}_{t-1}\right]\right)^2}{I_t(X^*, Y_t)} = \frac{(\text{expected regret})^2}{\text{mutual information}}$$

• Mutual information measures information gain

$$I_t(X^*, Y_t) = E\left[H_t(X^*) - H_{t+1}(X^*)|\mathbb{F}_{t-1}\right]$$

• Information ratio (IR)

$$\Psi_t = \frac{\left(E\left[f_{\theta}(X^*) - f_{\theta}(X_t) | \mathbb{F}_{t-1}\right]\right)^2}{I_t(X^*, Y_t)} = \frac{(\text{expected regret})^2}{\text{mutual information}}$$

• Mutual information measures information gain

$$I_t(X^*, Y_t) = E\left[H_t(X^*) - H_{t+1}(X^*) | \mathbb{F}_{t-1}\right]$$

• Entropy $H_t(X^*)$ measures degree of uncertainty

• Information ratio (IR)

$$\Psi_t = \frac{\left(E\left[f_{\theta}(X^*) - f_{\theta}(X_t) | \mathbb{F}_{t-1}\right]\right)^2}{I_t(X^*, Y_t)} = \frac{(\text{expected regret})^2}{\text{mutual information}}$$

• Mutual information measures information gain

$$I_t(X^*, Y_t) = E\left[H_t(X^*) - H_{t+1}(X^*) | \mathbb{F}_{t-1}\right]$$

- Entropy $H_t(X^*)$ measures degree of uncertainty
- IDS: select action distribution that minimizes Ψ_t

• Information ratio (IR)

$$\Psi_t = \frac{\left(E\left[f_{\theta}(X^*) - f_{\theta}(X_t) | \mathbb{F}_{t-1}\right]\right)^2}{I_t(X^*, Y_t)} = \frac{(\text{expected regret})^2}{\text{mutual information}}$$

• Mutual information measures information gain

$$I_t(X^*, Y_t) = E\left[H_t(X^*) - H_{t+1}(X^*) | \mathbb{F}_{t-1}\right]$$

- Entropy $H_t(X^*)$ measures degree of uncertainty
- IDS: select action distribution that minimizes Ψ_t
 - Trades off between expected regret and information gain

• Information ratio (IR)

$$\Psi_t = \frac{\left(E\left[f_{\theta}(X^*) - f_{\theta}(X_t) | \mathbb{F}_{t-1}\right]\right)^2}{I_t(X^*, Y_t)} = \frac{(\text{expected regret})^2}{\text{mutual information}}$$

• Mutual information measures information gain

$$I_t(X^*, Y_t) = E\left[H_t(X^*) - H_{t+1}(X^*) | \mathbb{F}_{t-1}\right]$$

- Entropy $H_t(X^*)$ measures degree of uncertainty
- IDS: select action distribution that minimizes Ψ_t
 - Trades off between expected regret and information gain
 - Support is of cardinality at most 2

Learning to Optimize

Relation to TS and Regret Bound

Relation to TS and Regret Bound

• A regret bound that applies to all algorithms

$$E[\operatorname{Regret}(T)] \le \sqrt{\overline{\Psi}_T H(X^*)T}$$

[Russo-Van Roy, 2014]

$$\overline{\Psi}_T = \frac{1}{T} \sum_{t=1}^T E[\Psi_t]$$
• A regret bound that applies to all algorithms

$$E[\operatorname{Regret}(T)] \le \sqrt{\overline{\Psi}_T H(X^*)T}$$

[Russo-Van Roy, 2014]

$$\overline{\Psi}_T = \frac{1}{T} \sum_{t=1}^T E[\Psi_t]$$

• For IDS:

• A regret bound that applies to all algorithms

$$E[\operatorname{Regret}(T)] \leq \sqrt{\overline{\Psi}_T H(X^*)T}$$

[Russo-Van Roy, 2014]

$$\overline{\Psi}_T = \frac{1}{T} \sum_{t=1}^T E[\Psi_t]$$

• For IDS: • $\overline{\Psi}_t \leq |\mathbb{X}|/2$ always

• A regret bound that applies to all algorithms

$$E[\operatorname{Regret}(T)] \leq \sqrt{\overline{\Psi}_T H(X^*)T}$$

[Russo-Van Roy, 2014]

$$\overline{\Psi}_T = \frac{1}{T} \sum_{t=1}^T E[\Psi_t]$$

- For IDS:
 - $\overline{\Psi}_t \leq |\mathbb{X}|/2$ always
 - $\overline{\Psi}_t \leq d/2$ for *d*-dimensional linear bandit

• A regret bound that applies to all algorithms

$$E[\operatorname{Regret}(T)] \leq \sqrt{\overline{\Psi}_T H(X^*)T}$$

[Russo-Van Roy, 2014]

$$\overline{\Psi}_T = \frac{1}{T} \sum_{t=1}^T E[\Psi_t]$$

- For IDS:
 - $\overline{\Psi}_t \leq |\mathbb{X}|/2$ always
 - $\overline{\Psi}_t \leq d/2$ for *d*-dimensional linear bandit
 - 1/2 with full feedback

• A regret bound that applies to all algorithms

$$E[\operatorname{Regret}(T)] \leq \sqrt{\overline{\Psi}_T H(X^*)T}$$

[Russo-Van Roy, 2014]

$$\overline{\Psi}_T = \frac{1}{T} \sum_{t=1}^T E[\Psi_t]$$

- For IDS:
 - $\overline{\Psi}_t \leq |\mathbb{X}|/2$ always
 - $\overline{\Psi}_t \leq d/2$ for *d*-dimensional linear bandit
 - 1/2 with full feedback
- Grew out of information-theoretic analysis of TS

[Russo-Van Roy, 2014]

Luntern 2015

Learning to Optimize

• Tractable implementations for several cases

- Tractable implementations for several cases
 - Beta-Bernouli bandit (independent arms)

- Tractable implementations for several cases
 - Beta-Bernouli bandit (independent arms)
 - Gaussian bandit (independent arms)

- Tractable implementations for several cases
 - Beta-Bernouli bandit (independent arms)
 - Gaussian bandit (independent arms)
 - Linear bandit (mean-based IDS)

- Tractable implementations for several cases
 - Beta-Bernouli bandit (independent arms)
 - Gaussian bandit (independent arms)
 - Linear bandit (mean-based IDS)
- UCB/TS do well in these cases

- Tractable implementations for several cases
 - Beta-Bernouli bandit (independent arms)
 - Gaussian bandit (independent arms)
 - Linear bandit (mean-based IDS)
- UCB/TS do well in these cases
- New algorithms needed for other cases

Linear Bandit Simulation



Linear Bandit Simulation



Learning to Optimize

• IDS addresses cases where UCB/TS miserably fails

- IDS addresses cases where UCB/TS miserably fails
- IDS accomplishes this by measuring information gain

- IDS addresses cases where UCB/TS miserably fails
- IDS accomplishes this by measuring information gain
- IDS performs as well or better than UCB/TS in several cases where all are tractable

- IDS addresses cases where UCB/TS miserably fails
- IDS accomplishes this by measuring information gain
- IDS performs as well or better than UCB/TS in several cases where all are tractable
- New algorithms are needed to implement IDS in other cases, especially those in which UCB/TS miserably fail