1) Model-Based RL in Query model
   (MBRL)

   Query model: any $(s,a)$ we can observe sample
   $$s' \sim P(s,a) \quad (\text{or } s' = f(s,a,w) \text{ w/ } w \sim D)$$

   Black-box
       - games
       - Simulator

   Sample complexity: How many samples are required
   for near-optimal performance?

Meta-Algorithm (MBRL)
   1) For $i = 1, \dots, N$        specify
      sample $s'_i \sim P(s_i, a_i)$   record $(s'_i, s_i, a_i)$

   2) Fit transition model $\hat{P}$ using $\{(s'_i, s_i, a_i)\}_{i=1}^{N}$
      specify
   3) Design $\hat{\pi}$ using $\hat{P}$
      specify

2) Tabular Setting                           $(N > SA)$
   1) samples $(s_i, a_i)$ evenly: $\frac{N}{SA}$ times each

   2) Fit transition model
      $$\hat{P}(s'|s,a) = \frac{\sum_{i=1}^{N} \mathbb{1}\{s_i = s \ \& \ a_i = a\} \mathbb{1}\{s'_i = s'\}}{\sum_{i=1}^{N} \mathbb{1}\{s_i = s, \ a_i = a\}}$$

   3) Design $\hat{\pi}$ Policy Iteration $\hat{\pi} = PI(\hat{P}, r)$

Recall: $PI(P,r)$

Initialize $\pi^0$

For $t=1,\dots,T$

$\qquad Q^{\pi^t} = $ Policy Eval$(\pi^t, P, r)$ $\quad\nearrow\quad \dfrac{V^{\pi^t}}{Q^{\pi^t}(s,a)} = \dfrac{(I-\gamma P^\pi)^{-1} R^\pi}{r(s,a) + \gamma \sum_{s'} V^\pi(s')}$

$\qquad \rightarrow \pi^{t+1}(s) = \underset{a \in \mathcal{A}}{\arg\max}\ \underline{Q^{\pi^t}(s,a)}\ \forall s$

$V^\pi(s) = Q^\pi(s, \pi(s))$

Goal: compare $\pi^*$ vs. $\hat{\pi}$

strategy: 
 I) compare $\hat{P}$ vs. $P$

 II) Translate $\hat{P}$ vs. $P$ to $\hat{V}$ vs. $V$

 III) translate $\hat{V}$ vs. $V$ to $\hat{\pi}$ vs. $\pi_*$

I) Model Estimation

 <u>Lemma</u>: With probability $1-\delta$, $\forall s,a$ $\quad SA - s$

$$\underbrace{\sum_{s' \in S} |\hat{P}(s'|s,a) - P(s'|s,a)|}_{\|\hat{P}(\cdot|s,a) - P(\cdot|s,a)\|_1} \leq \sqrt{\frac{S^2 A \log(2SA/\delta)}{\underline{N}}}$$

 Proof is out of scope

II) Value Functions.

$$V^\pi(s) = \mathbb{E}\left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \ \middle|\ \begin{matrix} s_0 = s \\ P \\ \pi \end{matrix} \right]$$

$$\hat{V}^\pi(s) = \mathbb{E}\left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \ \middle|\ \begin{matrix} s_0 = s \\ \hat{P} \\ \pi \end{matrix} \right]$$

Recall: Discounted State-Action Distribution

$$d_{s_0}^\pi(s,a) = (1-\gamma)\sum_{t=0}^{\infty}\gamma^t \underbrace{P_t^\pi(s,a;s_0)}_{\text{prob. of }(s,a)\text{ @ }t\text{ given }s_0\,\&\,\pi,P}$$

Simulation Lemma: $0 \le r(s,a) \le 1$

$$\|\hat{V}^\pi(s_0) - V^\pi(s_0)\| \le \frac{\gamma}{(1-\gamma)^2}\,\underset{s,a\sim d_{s_0}^\pi}{\mathbb{E}}\left[\|\hat{P}(\cdot|s,a) - P(\cdot|s,a)\|_1\right]$$

$\forall s_0$

under $P$ (arrow to $d_{s_0}^\pi$)

disagreement $\hat{P}$ & $P$ (arrow to $\|\cdot\|_1$ term)

Proof: claim:

$$\hat{V}^\pi(s_0) - V^\pi(s_0) = \boxed{\frac{\gamma}{1-\gamma}\,\underset{s,a\sim d_{s_0}^\pi}{\mathbb{E}}\left[\underset{s'\sim\hat{P}(s,a)}{\mathbb{E}}(\hat{V}^\pi(s')) - \underset{s'\sim P(s,a)}{\mathbb{E}}(\hat{V}^\pi(s'))\right]}$$

(underbrace to $\bigstar$)

$$\bigstar = \sum_{s'\sim\mathcal{S}}\left[\hat{P}(s'|s,a) - P(s'|s,a)\right]\hat{V}^\pi(s')$$

using $r \le 1$    $\hat{V}^\pi(s') \le \frac{1}{1-\gamma}$

$$\le \boxed{\frac{1}{1-\gamma}}\sum_{s'\sim\mathcal{S}}|\hat{P}(s'|s,a) - P(s'|s,a)|$$

III) Policy Iteration    $\hat{\pi} = PI(\hat{P}, r)$    $\hat{\pi}$ is optimal for $\hat{P} \to$ no approx error on PI

$$V^*(s_0) - V^{\hat{\pi}}(s_0) \le V^*(s_0) - \hat{V}^{\pi^*}(s_0) + \hat{V}^{\hat{\pi}}(s_0) - V^{\hat{\pi}}(s_0)$$

($V^{\pi^*}$ label above, $\ge 0$ label)

$\hat{\pi}$ is optimal on $\hat{P}$, $\hat{V}^{\hat{\pi}}(s) \ge \hat{V}^\pi(s) \ \forall \pi$

$$\le \frac{1}{(1-\gamma)^2}\left(\underset{s,a\sim d_{s_0}^{\pi^*}}{\mathbb{E}}\left[\|\hat{P}(\cdot|s,a) - P(\cdot|s,a)\|_1\right] + \underset{s,a\sim d_{s_0}^{\hat{\pi}}}{\mathbb{E}}\left[\|\hat{P}(\cdot|s,a) - P(s,a)\|_1\right]\right)$$

$$\longrightarrow \leq \frac{2}{(1-\gamma)^2} \sqrt{\frac{S^2 A \log(2SA/\delta)}{N}} = \varepsilon$$

**Theorem:**

$$N = \frac{4 S^2 A \log\left(\frac{2SA}{\delta}\right)}{\varepsilon^2}, \quad \text{then } V^*(s_0) - V^{\hat{\pi}}(s_0) \leq \varepsilon$$
$$\text{w.p} \geq 1 - \delta$$

3) LQR

$$\min \mathbb{E}\left[\sum s_t^\top Q s_t + a_t^\top R a_t \mid s_{t+1} = A s_t + B a_t + w_t\right]$$
$$w \sim \mathcal{N}(0, \sigma^2 I)$$

MBRL:

1) i.i.d. $s_i \sim \mathcal{N}(0, \sigma^2 I_{n_s})$ $a_i \sim \mathcal{N}(0, \sigma^2 I_{n_a})$

2) estimate

$$(\hat{A}, \hat{B}) = \text{argmin} \sum_{i=1}^{N} \|s_i' - A s_i - B a_i\|_2^2$$

3) $\hat{K} = LQR(\hat{A}, \hat{B}, Q, R)$