

# Probabilistic Artificial Intelligence summary

Michael Etienne Van Huffel, Andrea Ghirlanda, Dr. Linda Maria de Cave

Probabilistic Artificial Intelligence summary created by *michavan@student.ethz.ch*, *aghirlanda@student.ethz.ch* and *ldecave@student.ethz.ch*

This summary has been written based on the Lecture 5263-5210-00 L Probabilistic Artificial Intelligence by Prof. A. Krause (Fall 2023). There is no guarantee for completeness and/or correctness regarding the content of this summary. Use it at your own discretion. The language used in this cheat sheet contains many abbreviations, which were necessary to adhere to the two-page limit imposed by the lecturer.

**Bayes (iid  $\Leftrightarrow$  uncorrelated,  $S = -\log u$ )**  
 $P(X_{1:n}) = P(X_1)P(X_2|X_1)\cdots P(X_n|X_{1:n-1})$   
**Bayes:**  $P(X|Y) = \frac{P(X,Y)}{P(Y)} = \frac{P(Y|X)P(X)}{P(Y)}$   
 $\text{var}(X \pm Y) = \text{var}(X) + \text{var}(Y) \pm 2\text{cov}(X, Y)$   
 $\mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$   
**Gauss:**  $\frac{1}{\sqrt{(2\pi)^d|\Sigma|}} \exp(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu))$   
**Cond:**  $P(X_A|X_B = x_B) = \mathcal{N}(\mu_{A|B}, \Sigma_{A|B})$   
 $\mu_A + \Sigma_{AB}\Sigma_{BB}^{-1}(x_B - \mu_B), \Sigma_{AA} - \Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA}$   
 $H(\prod p_i) = \sum_i H(p_i); H(N(\mu, \Sigma)) = \frac{1}{2}\ln|2\pi e\Sigma|$   
 $H(p, q) = H(p) + H(q|p), H(p||q) = H(p) + KL(p||q)$   
**Convex:**  $g(x) \searrow \Leftrightarrow x_1, x_2 \in \mathbb{R}, \lambda \in [0, 1] : g''(x) > 0; g(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda g(x_1) + (1 - \lambda)g(x_2)$  **Egg:**  $g \prec : g(E[X]) \leq E[g(X)]$   
**Bayes:** *Prior:*  $p(\theta)$ , *Like:*  $p(y_{1:n}|x_{1:n}, \theta) = \prod_{i=1}^n p(y_i|x_i, \theta)$ , *Post:*  $p(\theta|x_{1:n}, y_{1:n}) = \frac{1}{Z}p(\theta) \prod_{i=1}^n p(y_i|x_i, \theta)$ ,  $Z = \int (*d\theta)$ , *Pred:*  $p(y^*|x^*, x_{1:n}, y_{1:n}) = \int p(y^*|x^*, \theta)p(\theta|x_{1:n}, y_{1:n})d\theta$   
**BLR (Gauss prior/noise, ass. same as Ridge)**  
**BayI**  $w_{ls} = w_{MLE}, w_{ridg} = w_{MAP}(\lambda = \sigma_n^2/\sigma_p^2)$   
Test  $x^*, f^* = w^T x^*, y^* = f^* + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma_n^2)$   
 $p(w) = \mathcal{N}(0, \sigma_p^2 \mathbf{I}), p(w|X, y) = \mathcal{N}(w; \bar{\mu}, \bar{\Sigma})$   
 $p(y_i|x_i, w, \sigma_n) = \mathcal{N}(y_i; w^T x_i, \sigma_n^2), \bar{\mu} = \sigma_n^{-2} \bar{\Sigma} X^T y$   
 $\bar{\Sigma} = (\sigma_n^{-2} X^T X + \sigma_p^{-2} I)^{-1}, p(f^*|y^*, X, y, x^*) = \mathcal{N}(x^{*T} \bar{\mu}, x^{*T} \bar{\Sigma} x^* + \sigma_n^2)$ , **epistem/ aleator (irr)**  
 $\text{var}[y^*|x^*] = \mathbb{E}_\theta[\text{var}_{y^*}[y^*|x^*, \theta]] + \text{var}_\theta[\mathbb{E}_{y^*}[y^*|x^*, \theta]]$   
**RecUpt:** giv. prior  $p(\theta)$ , obs  $y_{1:n}, p^{(t)}(\theta) = p(\theta|y_{t+1})$  post. aftr  $t$  obs,  $p^{(t+1)}(\theta) = p(\theta|y_{1:t+1}) = p^{(t)} \cdot p(y_{t+1}|\theta), X_{t+1}^T X_{t+1} = X_t^T X_t + x_{t+1} x_{t+1}^T$   
 $X_{t+1}^T y_{t+1} = X_t^T y_t + y_{t+1} x_{t+1}$  **Func View:** instead of  $w \sim \mathcal{N}(0, \sigma_p^2 I)$ ,  $f$  prior  $f|X \sim \mathcal{N}(\Phi \mathbb{E}[w], \Phi \text{var}[w] \Phi^T) = \mathcal{N}(0, \sigma_p^2 \Phi \Phi^T)$ ,  
 $k(x, x') = \sigma_p^2 \cdot \phi(x)^T \phi(x') = \text{cov}[f(x), f(x')]$   
**CV:**  $\lambda = \hat{\sigma}_n^2 / \hat{\sigma}_p^2, \hat{\sigma}_n^2$  (MSE), find  $\hat{\sigma}_p^2 = \hat{\sigma}_n^2 / \lambda$   
 $X_{t+1} \perp X_{1:t-1}, Y_{1:t-1}|X_t, Y_t \perp Y_{1:t-1}|X_{t-1}, Y_t \perp X_{1:t-1}|X_t$   
State  $X_t$  obs  $Y_t$  Prior  $P(X_1) \sim \mathcal{N}(\mu, \Sigma)$   
**Mot:**  $X_{t+1} = F X_t + \epsilon_t, \epsilon_t \sim \mathcal{N}(0, \Sigma_x)$  **Sens:**  $Y_t = H X_t + \eta_t, \eta_t \sim \mathcal{N}(0, \Sigma_y)$  **Kupdt 4 pred:**  $\mu_{t+1} = F \mu_t + K_{t+1}(y_{t+1} - H F \mu_t) \Sigma_{t+1} = (I - K_{t+1} H)(F \Sigma_t F^T + \Sigma_x) \mathbf{K}_{gain}: K_{t+1} = (F \Sigma_t F^T + \Sigma_x) H^T (H(F \Sigma_t F^T + \Sigma_x) H^T + \Sigma_y)^{-1}$   
**GP**  $f \sim GP(\mu(x), K(x))$  ( **$\infty$ -dim Gaussian**)  
 $\infty$  set of RVs  $X$  s.t.  $\forall A \subseteq X, A = \{x_1, \dots, x_m\}$  it holds  $Y_A = [Y_{x_1}, \dots, Y_{x_m}] \sim \mathcal{N}(\mu_A, K_{AA})$  w\  $K_{AA}^{(ij)} = k(x_i, x_j)$  and  $\mu_A^{(i)} = \mu(x_i)$  (join Gaus)  
**Kern:** *sqr d exp anal*,  $\infty$  dif *exp* cont nodif, *Mat*( $\nu$ )  $\nu$  mal dif,  $\nu = 1/2$  Lapl,  $\nu = \infty$  Gaus  
**Cov**  $k$ : symmetric, PSD, *statry*:  $k(x, x') = k(x - x')$ , *isotropic*:  $k(x, x') = k(\|x - x'\|_2)$ .  
**Pred**  $p(f) = GP(f; \mu(x), k(x, x'))$ , observe

$y_i = f(x_i) + \epsilon_i, \epsilon_i \sim \mathcal{N}(0, \sigma^2), A = \{x_1, \dots, x_m\}$ . Often  $\mu(x) = 0, p(f|x_{1:m}, y_{1:m}) = GP(f; \mu', k')$  with  $k_{x,A} = [k(x, x_1) \dots k(x, x_m)]$ ,  $\mu'(x) = \mu(x) + \mathbf{k}_{x,A}(\mathbf{K}_{AA} + \sigma^2 \mathbf{I})^{-1}(\mathbf{y}_A - \mu_A)$   $k'(x, x') = k(x, x') - \mathbf{k}_{x,A}(\mathbf{K}_{AA} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_{x',A}^T$ , *Pred*  
*post:* at  $x^*, p(y^*|x_{1:m}, y_{1:m}, x^*) = \mathcal{N}(\mu_n^*, \sigma_n^{*2})$ ,  $\mu_n^* = \mu'(x^*), \sigma_n^{*2} = \sigma^2 + k'(x^*, x^*)$ , **Samp:** 1) discrete set to sample  $\mathbf{f} = [f_1 \dots f_n]$ ,  $\mathbf{f} \sim \mathbf{K}^{1/2} \epsilon = \mathbf{L} \epsilon$  2)  $p(f_1 \dots f_n) = \prod_{i=1}^n p(f_i|f_{1:i-1})$ , sample  $f_n \sim p(f_n|f_{1:n-1})$  **ModelSel:** 1) MLE  $\hat{\theta} = \text{amax}_\theta p(y|X, \theta)$ ,  $\int p(y|X, f)p(f|\theta)df$ ,  $\exp\{-\frac{1}{2} \mathbf{y}^T \mathbf{K}_y^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}_y|\}$  2) place hyper-prior  $p(\theta)$ , MAP  $\hat{\theta} = \text{amax}_\theta p(y|x, \theta)p(\theta)$ , FullB  $p(y^*|x^*, x_{1:n}, y_{1:n}) = \int p(y^*|x^*, f)p(f|x, \theta)p(\theta)d\theta$   
**Cost:**  $n$  variables, req. lin system time  $\mathcal{O}(n^3)$ , space  $\mathcal{O}(n^2)$ , **BLR:**  $\mathcal{O}(nd^2)$ . **Fast GP:** 1) *GPU* 2) *Local:* distance decaying kernel (e.g. RBF), only condition on points  $x'$  where  $|k(x, x')| > \tau$  3)a)  $k$  approx:  $k(x, x') \approx \phi(x)^T \phi(x'), \phi \in \mathbb{R}^m$ , then BLR  $\mathcal{O}(nm^2 + m^3)$  3)b) stat *true position* *Bochner*, **Random**  $k(x, x') = \int_{\mathbb{R}^d} p(\omega) e^{i \omega^T (x - x')} d\omega = \mathbb{E}_{\omega, b} [z_{\omega, b}(x) z_{\omega, b}(x')]$   $\approx \frac{1}{m} \sum_{i=1}^m z_{\omega^{(i)}, b^{(i)}}(x) z_{\omega^{(i)}, b^{(i)}}(x')$ ,  $\omega \sim p(\omega), b \sim \mathcal{U}(0, 2\pi]$ ,  $z_{\omega, b}(x) = (2/D)^{1/2} \cos(\omega^T \bar{x} \pm b)$   $\sigma_{\omega} = 1, \sigma_b = 1, \sigma_{\omega} = 1, \sigma_b = 1$  draw samples  $\omega^{(i)}, b^{(i)}$ ,  $k(x, x') \approx \phi(x)^T \phi(x')$  with  $\phi_i(x) = \frac{1}{\sqrt{D}} z_{\omega^{(i)}, b^{(i)}}(x)$   
**Rahm**  $M \subset \mathbb{R}^d$  compact for RFFs  $z(x)$ ,  $\sigma_p^2 = \mathbb{E}[\omega^T \omega], \mathbb{P}[\sup_{x, x' \in M} |z(x) - z(x')| \geq \epsilon] \leq 28 (\frac{\sigma_p}{\text{diam}(M)})^2 \exp(-\frac{\epsilon^2}{4(d+2)})$   
4) *Inducing points methods:* Summarize data via values of  $f$  at inducing points  $\mathbf{u} = \{u_1, \dots, u_m\}$ .  
 $p(\mathbf{f}^*, \mathbf{f}) = \int p(\mathbf{f}^*, \mathbf{f}, \mathbf{u}) d\mathbf{u} = \int p(\mathbf{f}^*, \mathbf{f}|\mathbf{u}) p(\mathbf{u}) d\mathbf{u}$   
 $p(\mathbf{f}^*, \mathbf{f}) \approx q(\mathbf{f}^*, \mathbf{f}) = \int q(\mathbf{f}^*|\mathbf{u}) q(\mathbf{f}|\mathbf{u}) p(\mathbf{u}) d\mathbf{u}$  with  $u = f(z)$ ,  $z$ -inducing location, train conditional:  $p(\mathbf{f}|\mathbf{u}) = \mathcal{N}(\mathbf{K}_{f,u} \mathbf{K}_{u,u}^{-1} \mathbf{u}, \mathbf{K}_{f,f} - \mathbf{Q}_{f,f})$ , w\  $\mathbf{Q}_{a,b} = \mathbf{K}_{a,u} \mathbf{K}_{u,u}^{-1} \mathbf{K}_{u,b}$ , test cond:  $p(\mathbf{f}^*|\mathbf{u}) = \mathcal{N}(\mathbf{K}_{f^*,u} \mathbf{K}_{u,u}^{-1} \mathbf{u}, \mathbf{K}_{f^*,f^*} - \mathbf{Q}_{f^*,f^*})$   
4)a) *Subset of Regressors:* assume  $\mathbf{K}_{f,f} - \mathbf{Q}_{f,f} = 0$ , replace  $p(\mathbf{f}|\mathbf{u})$  by  $q_{SoR}(\mathbf{f}|\mathbf{u}) = \mathcal{N}(\mathbf{K}_{f,u} \mathbf{K}_{u,u}^{-1} \mathbf{u}, 0)$  resulting model is degenerate GP with covariance function  $k_{SoR}(\mathbf{x}, \mathbf{x}') = k(x, \mathbf{u}) \mathbf{K}_{u,u}^{-1} k(\mathbf{u}, \mathbf{x}')$  **FITC:** Assume  $\mathbf{f}_i \perp \mathbf{f}_j | \mathbf{u}, \forall i \neq j$   $q_{FITC}(\mathbf{f}|\mathbf{u}) = \mathcal{N}(\mathbf{K}_{f,u} \mathbf{K}_{u,u}^{-1} \mathbf{u}, \text{diag}(\mathbf{K}_{f,f} - \mathbf{Q}_{f,f}))$ ,  $q_{FITC}(\mathbf{f}^*|\mathbf{u}) = p(\mathbf{f}^*|\mathbf{u})$ , *cost* cubic in # inducing pts, dominated inv  $\mathbf{K}_{u,u}$ . *Pick inducing pts* by chose randly/greed criterion (var)/det grid,  $\mathbf{u}$  as hyperpar, max marg like (ensure  $\mathbf{u}$  repr data).

**Variational Inf (rev greed sel mode, fwd sel var)**  
 $\text{amin}_\lambda KL(p||q_\lambda) = \text{amax}_\lambda \lim_{n \rightarrow \infty} \sum_{i=1}^n \log q(x^{(i)}|\lambda)$   
 $q = \arg \min_{q \in \mathcal{Q}} KL(p||q)$  match 1, 2 mom of  $p$   
**Laplace:**  $p(\theta|(x, y)_{1:n}) \approx q_\lambda(\theta) = \mathcal{N}(\hat{\theta}, \Lambda^{-1})$   
 $\hat{\theta} = \arg \max_\theta p(\theta|y), \Lambda = -\nabla^2 \log p(\hat{\theta}|y), 0 = \nabla \log p(\hat{\theta}|y)$ , *Pred* BL, greed fit mode (overcnf), match curv, pres MAP est, can diff from post  
**GVI:**  $0 = \mathbb{E}_{q(\theta)} [\nabla_\theta \log p(\mathcal{D}, \theta)], \Sigma^{-1} = -\mathbb{E}_{q(\theta)} [\nabla_\theta^2 \log p(\mathcal{D}, \theta)]$  **PVI:**  $KL(\mathcal{N}||Po) \propto -d/2 \ln(2\pi e) - d \ln \sigma - 1/m \sum_{j,i=1}^{m,n} [y_i w^T x_i - e^{w^T x_i - \|w\|_2^2 / 2\sigma_p^2}]$  **BLR:**  $P(y|x, w) = \text{Ber}(y; \sigma(w^T x))$ ,  $\hat{w} = \text{amin}_w \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \lambda \|w\|_2^2$ ,  $\mathcal{N}(\hat{w}, \Lambda^{-1}), \Lambda = \sum_{i=1}^n x_i x_i^T \sigma(\hat{w}^T x_i)(1 - \sigma(\hat{w}^T x_i))$   
**Variational Inf:**  $p(\theta|y) = \frac{1}{Z} p(\theta, y) \approx q_\lambda(\theta)$   
 $q^* \in \arg \min_{q \in \mathcal{Q}} KL(q||p): q \approx p$  where  $\mathbf{q}$  large  $\text{amin}_q KL(q||p) = \text{amax}_q \mathbb{E}_{\theta \sim q} [\log p(\theta, y)] + H(q(\theta)) = \text{amax}_q \mathbb{E}_{\theta \sim q_\lambda(\theta)} [\log p(y|\theta)] - KL(q(\theta)||p(\theta))$   
 $\max_\epsilon \mathbb{E}_\epsilon [\ln p(\mathcal{D}, \mu + \Sigma^{1/2} \epsilon)] + \frac{1}{2} \ln |\Sigma| + \frac{d}{2} \log 2\pi e$   
**ELBO,  $L(\lambda) \leq \log p(y)$ .**  $\nabla_\lambda L(\lambda)$  hard, *score*  $\nabla: \nabla_\lambda L = \mathbb{E}_{\theta \sim q_\lambda} [\nabla_\lambda \log q(\theta|\lambda)(\log p(y, \theta) - \log q(\theta|\lambda))]$   
 $\theta \sim q_\lambda(\cdot)$  dep on var param. **Reparam.:** let  $\epsilon \sim \phi, \theta = g(\epsilon, \lambda), q(\theta|\lambda) = \phi(\epsilon) |\nabla_\epsilon g(\epsilon; \lambda)|^{-1}$  and  $\mathbb{E}_{\theta \sim q_\lambda} [f(\theta)] = \mathbb{E}_{\epsilon \sim \phi} [f(g(\epsilon; \lambda))]$ , which yield  $\nabla_\lambda \mathbb{E}_{\theta \sim q_\lambda} [f(\theta)] = \mathbb{E}_{\epsilon \sim \phi} [\nabla_\lambda f(g(\epsilon; \lambda))]$   
**Blackbox VI:** max ELBO using stoch opt., for diagonal  $q$ , 2x expensive as MAP, only need to diff joint prob  $p$  and  $q$ , also use Natural Grad, Variance reduct tech. **GP Class:**  $P(f) = GP(\mu, k), P(y|f, \mathbf{x}) = \sigma(y \cdot f(\mathbf{x}))$  max (\*) w\  $q(f_i) := \int p(f_i|\mathbf{u}) q(\mathbf{u}) d\mathbf{u}, \mathbf{u}$  pseudo in.  
 $\sum_{i=1}^n \mathbb{E}_{q(f_i)} [\log p(y_i|f_i)] - KL(q(\mathbf{u})||p(\mathbf{u}))$   
**MCMC (MC: seq of RV  $X_{1:n}$  with \*)**  
**TV Dist:**  $\|\mu - v\|_{\text{TV}} = 2 \sup_{A \subseteq \mathcal{A}} |\mu(A) - v(A)|$   
**Mix.Time:**  $\tau_{\text{TV}}(\epsilon) = \min\{t | \forall q_0: \|q_t - \pi\|_{\text{TV}} \leq \epsilon\}$   
*Rapidly mix:*  $\tau_{\text{TV}}(\epsilon) \in \mathcal{O}(\text{poly}(n, \log(1/\epsilon)))$   
(\*) Stat MC w\ prior  $P(X_1)$ , trans  $P(X_{t+1}|X_t)$  indep of  $t$ . *ergodic*  $\exists t < \infty$  s.t. all states reachable from every state in *exactly*  $t$  steps. **Mark Ass:**  $X_{t+1} \perp\!\!\!\perp X_{1:t-1} | X_t \forall t$  **Stat Distrib:**  $\pi(P - I) = 0$ . Stat Ergodic MC is unique stat distr  $\pi(X) > 0$  s.t.  $\forall x: \lim_{N \rightarrow \infty} P(X_N = x) = \pi(x), \pi(X)$  indep of prior  $P(X_1)$ . **Sim MC** fwd sampl  $x_N \sim P(X_N|X_{N-1} = x_{N-1})$  **MCMC:** Approx pred. distr.  $p(y^*|x^*, x_{1:n}, y_{1:n}) = \int p(y^*|x^*, \theta)p(\theta|(x, y)_{1:n})d\theta = \mathbb{E}_{\theta \sim p(\cdot|(x, y)_{1:n})} [p(y^*|x^*, \theta)] \approx \frac{1}{m} \sum_{i=1}^m p(y^*|x^*, \theta^{(i)})$ , sample  $\theta^{(i)} \sim p(\theta|(x, y)_{1:n})$  from MC with stat distr (\*).  
**Hoeffding:** Let  $f$  be bounded in  $[0, C]: \mathbb{P}(|\mathbb{E}[f(X)] - \frac{1}{N} \sum_{i=1}^N f(x_i)| > \epsilon) \leq 2 \exp(-2N\epsilon^2/C^2)$   
Given unnormalized distr.  $Q(x) > 0$ , design MC s.t.  $\pi(x) = \frac{1}{Z} Q(x)$ . **DB (reversible):**  $Q(x)P(x'|x) = Q(x')P(x|x') \rightarrow \pi(x) = \frac{1}{Z} Q(x)$ .

**Gibbs Sampling:** Asympt. correct but slow  
1. Init  $\mathbf{x}^{(0)}$ , fix observed RVs  $X_B$  to  $\mathbf{x}_B$   
2. Repeat: **set**  $\mathbf{x}^{(t)} = \mathbf{x}^{(t-1)}$ ; **sel**  $j \in [1:m] \setminus B$   $x_j^{(t)} \sim P(X_j|\mathbf{x}_{[1:m] \setminus \{j\}}^{(t)})$  **Rand:** ful. DB, find cor distr. **Determin:** not ful. DB, corr distr.  $P(X_i = x_i | \mathbf{x}_{-i}) = \frac{1}{Z} Q(X_i = x_i, \mathbf{x}_{-i}) = \frac{1}{Z} Q(\mathbf{x}_{1:n})$  re-sampl  $X_i$  only requires eval unnorm joint distrib and renorm. **Expectations via MCMC:** Joint sample at  $t$  dep only on  $t-1 \rightarrow$  LLN, HB not apply. *Thm:*  $X_{1:n}$  EMC on finite  $D, f \in D$ ,  $\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N f(x_i) = \sum_{x \in D} \pi(x) f(x)$   
Use MCMC to get samples  $\mathbf{X}^{(1:T)}$ . After burn-in time  $t_0: \mathbb{E}[f(\mathbf{X})] \approx \frac{1}{T-t_0} \sum_{\tau=t_0+1}^T f(\mathbf{X}^{(\tau)})$   
**MetroH:** Gen MC  $\rightarrow$  DB 1) Proposal  $R(X'|X)$ , given  $X_t = x$ , sample  $x' \sim R(X'|X = x)$ ; 2) For  $X_t = x$ , w.p.  $\alpha = \min\{1, \frac{Q(x')R(x|x')}{Q(x)R(x'|x)}\}$ :  $X_{t+1} = x'$ ; else  $X_{t+1} = x$  **Cont RVs:** log-concave  $p(x) = \frac{1}{Z} \exp(-f(x))$ ,  $f$  convex. M/H:  $\alpha = \min\{1, \frac{R(x|x')}{R(x'|x)} \exp(f(x) - f(x'))\}$ , Simple uninfl dir. **Improved prop: MALA/LMC:**  $R(x'|x) = \mathcal{N}(x'; x - \eta_t \nabla f(x); 2\tau I)$  MALA converg to stat distr for log concave (locall also non convex) func (for general distrib convergence slow), mixing time  $\mathcal{O}(d)$ . *Improve efficiency:* both proposal step and accept step requires full acces to energy function  $f \rightarrow$  SGD, decaying step size, skip accept/reject  $\rightarrow$  SGLD  
 $\theta \sim \frac{1}{Z} \exp(\log p(\theta) + \sum_{i=1}^n \log p(y_i|x_i, \theta))$  1)  $\theta_0$  2) For  $t$  do  $i_{1..i_m} \sim \mathcal{U}(1, n), \epsilon_t \sim \mathcal{N}(0, 2\eta_t I)$ ,  $\theta_{t+1} = \theta_t - \eta_t (\nabla \ln p_{\theta_t} + \frac{n}{m} \sum_{i=1}^m \nabla \ln p(y_{i_j}|\theta_t, x_{i_j})) + \epsilon_t$  SGLD = SGD + Gauss noise, convergence if  $\eta_t \in \mathcal{O}(t^{-1/3})$ , const step boost mixing, improve perf. via Adagrad, *HMC* (mom)  
**BDL (Prior:  $p(\theta) = \mathcal{N}(\theta; 0, \sigma_p^2 I)$ , Gauss = \*)**  
\* weig decay, *Like:*  $p(y|\mathbf{x}, \theta) = \mathcal{N}(y; f(\mathbf{x}, \theta), \sigma^2)$   
MAP:  $\hat{\theta} = \text{amin}_\theta -\log p(\theta) - \sum_i \log p(y_i|x_i, \theta)$  hetero  $\epsilon$  well, fails pred epistemic, use VI (*BbB*). **Etero:**  $y|\mathbf{x}, \theta \sim \mathcal{N}(\mu(\mathbf{x}; \theta), \sigma^2(\mathbf{x}; \theta))$ ,  $\mathbf{f}_1(\mathbf{x}; \theta), \exp(\mathbf{f}_2(\mathbf{x}; \theta))$  **VI:** SGD-opt ELBO via  $\nabla_\lambda L(\lambda)$ . Find VI approx  $q_\lambda$ . Draw  $m$  weights  $\theta^{(j)} \sim q_\lambda(\cdot)$ . *Pred*  $p(y^*|\mathbf{x}^*, \mathbf{x}_{1:n}, y_{1:n}) \approx \mathbb{E}_{\theta \sim q_\lambda} [p(y^*|\mathbf{x}^*, \theta)] \approx \frac{1}{m} \sum_j p(y^*|\mathbf{x}^*, \theta^{(j)})$ ,  $\text{Var}[y^*|\cdot] \approx \frac{1}{m} \sum_1^m \sigma^2(x^*, \theta^{(j)}) + \frac{1}{m} \sum_1^m (\mu(x^*, \theta^{(j)}) - \bar{\mu}(x^*))^2$   
**MC** wghts  $\theta^{(1:T)}$  SGLD, LD, SGHMC; pred by avg  $\theta^{(1:T)}$ . *Summ:* subsmpl/GVI  $q(\theta|\mu_{1:d}, \sigma_{1:d}^2)$   
 $\mu_i = \frac{1}{T} \sum_{j=1}^T \theta_i^{(j)}; \sigma_i^2 = \frac{1}{T} \sum_{j=1}^T (\theta_i^{(j)} - \mu_i)^2$   
**DropVI:**  $q_j(\theta_j|\lambda_j) = p\delta_0(\theta_j) + (1-p)\delta_{\lambda_j}(\theta_j), \theta^{(j)}$   
NN w\ wght giv by  $\lambda, \theta^{(j)} = 0$  wp  $p$ .  
**Prob. Ensbles:** train multiple models bootstrap, **Cal:** confidence=accuracy on heldout  
**Reliability Diag:** plot expected sample acc



as func of confidence 1) group pred in  $M$  bin, 2)  $\text{freq}(B_m) = 1/|B_m| \sum_{i \in B_m} [\hat{y}_i = 1]$  3)  $\text{conf}(B_m) = 1/|B_m| \sum_{i \in B_m} \hat{p}_i$  4) **MACE** =  $\max \sum_{m=1}^M \frac{|B_m|}{n} |\text{freq}(B_m) - \text{conf}(B_m)|$  Improve via histo bin, isotonic regr., platt (temp) scaling

**Active Learn (min #x reducing uncertainty)**

**MutInf:**  $I(X;Y) = I(Y;X) = H(X) - H(X|Y)$   
 $X \sim N(\mu, \Sigma), Y = X + N(0, \sigma^2), I = \frac{1}{2} \ln |I + \frac{1}{\sigma^2} \Sigma|$

**InfGain:**  $f(S), S \subseteq D, F(S) = H(f) - H(f|y_S) = I(f; y_S) = \frac{1}{2} \log |I + \sigma^{-2} K_S|$ , obs at  $S$   
*Thm:*  $F(S_T) \geq (1 - 1/e) \max_{S \subseteq D, |S| \leq T} F(S)$

**Greedy MIOpt:**  $F(S)$  NP hard,  $S_t = \{x_1, \dots, x_t\}$   
 $x_{t+1} = \arg \max_{x \in D} F(S_t \cup \{x\}) = \arg \max_{x \in D} \sigma_{x|S_t}^2$

**UncSmpl:**  $x_t = \arg \max_{x \in D} \sigma_{t-1}^2(x)$ , (homo,  $\mathcal{N}$ )  
 Fail dist **epist/aleat**, *Etero:*  $\text{amax}_x \sigma_f^2(x) / \sigma_n^2(x)$

*Or:* trAce/Eigval/log-Determinant designs

**BALD:** (class)  $x_{t+1} = \text{amax}_x I(\theta; y_x | x_{1:t}, y_{1:t}) = \text{amax}_x H(y|x, (x, y)_{1:t}) - \mathbb{E}_{\theta \sim p(\cdot | (x, y)_{1:t})} [H(y|x, \theta)]$

**Bayesian Opt (seq. pick  $x_1, \dots, x_T \in D$  \*)**  
 (\*),  $y_t = f(x_t) + \epsilon_t$ , find  $\max_x f(x)$  st  $T$  smal

**CumReg:**  $R_T = \sum_{t=1}^T \max_{x \in D} f(x) - f(x_t)$   
 sublin if  $R_T/T \rightarrow 0$  (obj)  $\Leftrightarrow f(x_t) \rightarrow \max f(x)$   
 (UCB  $\geq$  best lower bound if well cal) **GP-UCB:**  
 $x_t = \arg \max_{x \in D} \mu_{t-1}(x) + \beta_t \sigma_{t-1}(x)$   
 $\mu(x), \sigma(x)$  from GP marginal,  $\beta_t$  EE-tradeoff. non-convex usually, low  $D$  use Lipschitz, high  $D$ . grad ascent (w\ rand init) *Thm:*  $f \sim GP$ , if correct  $\beta_t: \frac{1}{T} R_T = \mathcal{O}(\sqrt{\gamma_T/T})$ , with  $\gamma_T = \max_{|S| \leq T} I(f; y_S)$  (max. info gain)

**Thomp sample:** at  $t$ , draw from GP post.  $\tilde{f} \sim P(f|x_{1:t}, y_{1:t})$ , select  $x_{t+1} \in \arg \max_{x \in D} \tilde{f}(x)$   
 datasets in BO/AL small, data pts selected dependend on prior obs, *sol:* hyperprior on hyperparam, select pts at random occasionally.

**Markov Decision Processes ( $r, P$  known)**

**MDP:** Finite MDP (control MC), state  $X = \{1, \dots, n\}$ , action  $A = \{1, \dots, m\}$ , trans  $P(x'|x, a)$ , init  $P(x_0)$ , reward  $r(x, a, x')$ , discount  $\gamma \in [0, 1]$

**Planning in MDPs:** Policy  $\pi: X \rightarrow AP(A)$  (**det./rand**), induce MC w\ trans  $P(X_{t+1} = x' | X_t = x) = P(x' | x, \pi(x)) \sum_a \pi(a | x) P(x' | x, a)$

**Value fun:** fixed  $\pi$   $V^\pi(x) = J(\pi | X_0 = x) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r(X_t, \pi(X_t)) | X_0 = x] = r(x, \pi(x)) + \gamma \sum_{x'} P(x' | x, \pi(x)) V^\pi(x')$ ,  $V^\pi = (I - \gamma T^\pi)^{-1} r^\pi$   
 $V_i^\pi = V^\pi(i)$ ,  $r_i^\pi = r^\pi(i, \pi(i))$ ,  $T_{i,j}^\pi = P(j | i, \pi(i))$   
 $V^\pi(x) = \sum_{x'} P(x' | x, \pi(x)) [r(x, \pi(x), x') + \gamma V^\pi(x')]$   
 $= Q^\pi(x, \pi(x)) = \mathbb{E}_{a' \sim \pi(x)} Q^\pi(x, a')$  **det./rand**

**Fixed Point Iteration:** 1) initialize  $V_0^\pi$  2) for  $t = 1 : T$  do:  $V_t^\pi = r^\pi + \gamma T^\pi V_{t-1}^\pi$  (converges)

**Greedy policy w.r.t.  $V$ :**  $V$  induces policy  $\pi_V(x) = \text{amax}_a r(x, a) + \gamma \sum_{x'} P(x' | x, a) V(x')$

**Greedy  $\pi$  w.r.t.  $Q$ :**  $\pi_Q(x) = \arg \max_a Q(x, a)$

**Bellman:** Opt  $\pi^* \leftrightarrow$  greedy wrt induced  $V^*$   
 $V^*(x) = \max_{a \in A} [r(x, a) + \gamma \sum_{x' \in X} P(x' | x, a) V^*(x')]$   
 $= \max_{a \in A} \mathbb{E}_{x'} [r(x, a) + \gamma V^*(x')] = \max_{a \in A} Q^*(x, a)$

**Policy Iter:** 1) Init arbitry  $\pi$  2) Until conv: calc  $V^{\pi_t}(x)$ , calc greedy pol  $\pi_t^G$  w.r.t.  $V^{\pi_t}$ , set  $\pi_{t+1} \leftarrow \pi_t^G$  Stop if  $V^{\pi_t}(x) = V^{\pi_{t+1}}(x)$ . Monoton improves all val  $V^{\pi_{t+1}}(x) \geq V^{\pi_t}(x) \forall x$ . Converge exact optimal in  $\mathcal{O}(n^2 m / (1 - \gamma))$ .

**SAF:**  $Q_t(x, a) = r(x, a) + \gamma \sum_{x'} P(x' | x, a) V_{t-1}(x')$

**Value Iteration:** 1) Init  $V_0(x) = \max_a r(x, a)$  2) for  $t = 1 : \infty$ :  $V_t(x) = \max_a Q_t(x, a)$ . Stop if  $\|V_t - V_{t-1}\|_\infty \leq \epsilon$ , pick grdy  $\pi_G$  w.r.t.  $V_t$ . Conv  $\epsilon$ -opt sol in poly time. **Tradeoffs:** **Pol/val**  $\mathcal{O}(n^3/nmk)$  per iter,  $k$  sprs **POMDP:**  $(X, A, P, R, \gamma, Y, Q)$ ,  $O_{x,y,a} = \frac{1}{|Y|} \sum_{y \in Y} [Y_{t+1} = y | X_{t+1} = x, A_t = a]$

**Belief state:**  $b_t(x) = P[X_t = x | y_{1:t}, A_{t-1} \equiv a]$   
 $b_t(x) \in \Delta^{|X|} = \{b \geq 0 \in \mathbb{R}^{|X|}, \sum_{i=1}^{|X|} b_i = 1\}$   
 $b_{t+1}(x) = \frac{1}{2} O(y_{t+1}, x) \Sigma_x b_t(x') P(x' | x, a)$ ,  $Z = \sum_x b_{t+1}(x)$

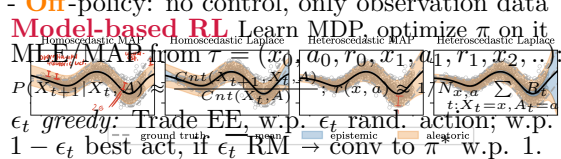
**BMDP:**  $(\Delta^{|X|}, A, \tau(b' | b, a), \rho = \sum_x b(x) r(x, a))$

**Solve POMDPs** finite  $T$ , exp. #belif states. BUT: most belif states never reach  $\rightarrow$  discretize space by sampling, PBVI, PBPI, dim reduction

**RL (agent act. chang state, unknown MDP)**

- **On-policy:** full control on actions/EE-trade  
 - **Off-policy:** no control, only observation data

**Model-based RL** Learn MDP, optimize  $\pi$  on it. **MEE/MAP** from  $\tau = (x_0, a_0, r_0, x_1, a_1, r_1, x_2, \dots)$



$\epsilon_t$  **greedy:** Trade EE, w.p.  $1 - \epsilon_t$  rand. action; w.p.  $\epsilon_t$  best act, if  $\epsilon_t$  RM  $\rightarrow$  conv to  $\pi^*$  w.p. 1.

**RobMonro (RM):**  $\sum_t \epsilon_t = \infty, \sum_t \epsilon_t^2 < \infty$

**$R_{\max}$  Alg:** Set unknown  $r(x, a)$  to  $R_{\max}$ ,  $r(x, a) \leq R_{\max}, \forall x, a$ , add **fairy tale state**  $x^*$ , set  $P(x^* | x, a) = 1$ , compute  $\pi$ . Repeat: run  $\pi$  while updtng  $r(x, a), P(x' | x, a)$ , recompute  $\pi$ .  $P(x^* | x^*, a) = 1, r(x^*, a) = R_{\max}$  *Lem:* evry  $T$  stps, whp  $R_{\max}$  eithr obt near opt reward, or visit at least 1 unknwn state-action pair. *Thm:* wp  $1 - \delta$ ,  $R_{\max}$  reach  $\epsilon$ -opt policy in #steps poly in  $|X|, |A|, T, 1/\epsilon, \log(1/\delta), R_{\max}$ .

**Problems of MBRL:** *Mem require:*  $P(x' | x, a) \approx \mathcal{O}(|X|^2 |A|)$ ,  $r(x, a) \approx \mathcal{O}(|X| |A|)$  *Computeate:* reptdly solve MDP (Val/Pol Iter)

**Model-free RL** Directly estim value function

**TD-Learning:** Follow policy  $\pi$ , get  $(x, a, r, x')$   
 Update:  $\hat{V}^\pi(x) \leftarrow (1 - \alpha_t) \hat{V}^\pi(x) + \alpha_t (r + \gamma \hat{V}^\pi(x'))$   
*Thm:* If  $\alpha_t$  is RM and all  $(x, a)$  pairs chosen  $\infty$  often, then  $\hat{V}^\pi$  converges to  $V^\pi$  w.p. 1.

**Optimistic Q-L** Estimate  $Q^*(x, a)$  1) Init rand /zero  $\hat{Q}^*(x, a) = \frac{R_{\max}}{1 - \gamma} \prod_{i=1}^{T_{\text{init}}} (1 - \alpha_i)^{-1}$  2) at  $t$   $a_t \in \text{amax}_a \hat{Q}^*(x_t, a)$ , get  $(x_t, a_t, r, x')$ ,  $\hat{Q}^*(x_t, a_t) \leftarrow (1 - \alpha_t) \hat{Q}^*(x_t, a_t) + \alpha_t (r + \gamma \max_{a'} \hat{Q}^*(x', a'))$   
*Thms* for  $R_{\max}$ , TD-L  $\rightarrow$  Time:  $\mathcal{O}(|A|)$ , Mem:  $\mathcal{O}(|X| |A|)$ ,  $|X|, |A|$ , exp #agents, state vars

**RL via FuncApprox (parametric approx. of \*)**  
 (\*) val funct  $V(x; \theta)$  or act val funct  $Q(x, a; \theta)$

**TD-Learning as SGD:** Tabular TD update rule can be viewed as SGD on squared loss  $l_2(\theta; x, x', r) = \frac{1}{2} (V(x; \theta) - r - \gamma V(x'; \theta_{old}))^2$ , then  $V \leftarrow V - \alpha_t \nabla_{V(x; \theta)} l_2$  is equiv to TD update.

**Function Approx Q-learning:** very slow Loss  $l_2(\theta; x, a, r, x') = \frac{1}{2} \delta^2$  with  $\delta = Q(x, a; \theta) - r - \gamma \max_{a'} Q(x', a'; \theta_{old})$ . *Alg:* Until conv: In state  $x$ , pick action  $a$ , obsrv  $r, x'$ . Update:  $\theta \leftarrow \theta - \alpha_t \nabla_\theta l_2 = \theta - \alpha_t \delta \nabla_\theta Q(x, a; \theta)$

**DQN:** Q-learn w\ NN as func approx. *experience replay*, maintn data  $D$ , clone NN to stabilize target opt  $L(\theta) = \sum_{(x, a, r, x') \in D} (r + \gamma \max_{a'} Q(x', a'; \theta_{old}) - Q(x, a; \theta))^2$

**Double DQN:** Use current NN to eval action arg max to prevents maximiz. bias of DQN.  
 $L^{DDQN}(\theta) = \sum_{(x, a, r, x') \in D} [r + \gamma Q(x', a^*(\theta); \theta_{old}) - Q(x, a; \theta)]^2$   $a^*(\theta) = \arg \max_{a'} Q(x', a'; \theta)$

**Q-learn  $\pi$   $a_t = \text{amax}_a Q(x, a; \theta)$**  bad if  $|A| \uparrow$

**Policy Search Methods:** Param. policy  $\pi_\theta$  Max  $J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta(a|x)}$  ( $\tau = (x_0, a_0, r_0, x_1, a_1, r_1, \dots)$ )  
 $r(\tau) = \sum_{t=0}^T \gamma^t r(x_t, a_t)$ ; via  $\nabla_\theta$ . *Score Grad:*  $\nabla_\theta J_\theta = \nabla_\theta \mathbb{E}_{\tau \sim \pi_\theta} r(\tau) = \mathbb{E}_{\tau \sim \pi_\theta} [r(\tau) \nabla_\theta \ln \pi_\theta(\tau)]$   
 MDP:  $\pi_\theta(\tau) = p(x_0) \prod_0^T \pi(a_t | x_t; \theta) p(x_{t+1} | x_t, a_t)$  then:  $\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [r(\tau) \sum_{t=0}^T \nabla_\theta \log \pi(a_t | x_t; \theta)]$   
*Reduce variance:* **baselines** (but  $\nabla$  unbiased)  $\mathbb{E}_{\tau \sim \pi_\theta} [r(\tau) \nabla \log \pi_\theta(\tau)] = \mathbb{E} \dots [(r(\tau) - b) \nabla \log \pi_\theta(\tau)]$

**R2Go:**  $G_t = \sum_{t'=t}^T \gamma^{t'-t} r_{t'}$ ;  $b_t(x_t) = 1/T \sum_{t=0}^T G_t$   
 $\nabla J_T(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [\sum_{t=0}^T \gamma^t G_t \nabla_\theta \log \pi(a_t | x_t; \theta)]$   
*Mean over returns:* replace  $G_t$  with  $(G_t - b_t(x_t))$

**REINFORCE:** Input  $\pi(a|x; \theta)$ , init  $\theta$ , repeat: gener episode  $(x_i, a_i, r_i)$ ,  $i = 0 : T$ ; for  $t = 0 : T$  set  $G_t$  to retrn from step  $t$ , update  $\theta$ :  
 $\theta = \theta + \eta \gamma^t G_t \nabla_\theta \log \pi(a_t | x_t; \theta)$  optimizes score- $\nabla$  using MC returns; high variance

**Deep RL with policy grad and actor-critic**

**Advantage Func:**  $A^\pi(x, a) = Q^\pi(x, a) - V^\pi(x)$   
 $\pi^*$  iff  $\forall x, a : A^{\pi^*}(x, a) \leq 0; \forall \pi, x : \max_a A^\pi(x, a) \geq 0$

**Actor-Critic:** Approx  $\pi_\theta$  and  $V^\pi$ , e.g. 2 NNs  
 Reinterpreting score- $\nabla =$  *Policy grad Thm*  $\nabla J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [\sum_{t=0}^{\infty} \gamma^t Q(x_t, a_t; \theta_Q) \nabla \log \pi(a_t | x_t; \pi_\theta)]$   
 $= \mathbb{E}_{(x, a) \sim \pi_\theta} [Q(x, a; \theta_Q) \nabla_\theta \log \pi(a | x; \pi_\theta)]$

**Online A.C.:** val fun appr + pol grad thm  $\delta = Q(x, a; \theta_Q) - r - \gamma Q(x', \pi(x', \pi_\theta); \theta_Q)$   
 $\theta_\pi \leftarrow \theta_\pi + \eta_t Q(x, a; \theta_Q) \nabla \log \pi(a | x; \pi_\theta)$   
 $\theta_Q \leftarrow \theta_Q - \eta_t \delta \nabla Q(x, a; \theta_Q)$  (**FA Q-learning**)  
*Var red:* **baselines**  $Q(x, a; \theta_Q) - V(x; \theta_V)$ : adv func estim  $\rightarrow$  **A2/3C, GAE/GAAC, TRPO** opt surrogate obj in trust region, guarantes monoton improve  $J(\theta)$ , **PPO** heuristic variant

**Off-policy AC:** allow reuse of past data

**Replace in  $L_{DQN}$   $a'$  by  $\pi(x'; \pi_\theta)$** ,  $\pi$  follows greedy pol., to model  $\max_{a'} Q$ . Equivalent to:  $\theta_\pi^* \in \arg \max_\theta \mathbb{E}_{x \sim \mu} [Q(x, \pi(x; \theta); \theta_Q)]$ , with

$\mu(x) > 0$  explores all states. If  $Q(\cdot; \theta_Q), \pi(\cdot; \theta_\pi)$  diff, use backprop to get stoch- $\nabla$  (unbiased)  
 $\nabla_\theta J(\theta) = \mathbb{E}_{x \sim \mu} [\nabla_\theta Q(x, \pi(x; \theta); \theta_Q)]$   
 $\nabla_\theta Q(x, \pi(x; \theta)) = \nabla_a Q(x, a)|_{a=\pi(x; \theta)} \nabla_\theta \pi(x; \theta)$   
 Needs *deterministic*  $\pi \rightarrow$  inject additional action noise ( $\epsilon_t$  greedy) to ensure explore  $\rightarrow$

**DDPG:** 1) init  $\theta_Q, \theta_\pi$ , set  $\theta_Q^{old} = \theta_Q, \theta_\pi^{old} = \theta_\pi$  2) repeat: observe  $x$ , execute  $a = \pi(x; \theta_\pi) + \epsilon$  to observe  $r, x'$ , store in  $D$ . If time to update: for iter: sample batch  $B$  from  $D$ , cmpt target  $y = r + \gamma Q(x', \pi(x', \theta_\pi^{old}), \theta_Q^{old})$ , updates: *Critic:*  $\theta_Q \leftarrow \theta_Q - \eta \nabla 1/|B| \sum_B (Q(x, a; \theta_Q) - y)^2$ , *Actor:*  $\theta_\pi \leftarrow \theta_\pi + \eta \nabla 1/|B| \sum_B (Q(x, \pi(x; \theta_\pi); \theta_Q) - y)^2$   
*Params:*  $\theta_j^{old} \leftarrow (1 - \rho) \theta_j^{old} + \rho \theta_j$  for  $j \in \{\pi, Q\}$

**TD3:** DDPG with 2 Critics avoid maxim bias

**Rand  $\pi$  DDPG:** (DQN with reparam  $\pi$  grad)  
 For critic  $a' \sim \pi(x'; \theta_\pi^{old})$ , get unbias  $\nabla$ , *actor*  $\nabla_{\theta_\pi} \mathbb{E}_{a \sim \pi(x; \theta_\pi)} Q(x, a; \theta_Q)$  reparam  $a = \psi(x; \theta_\pi, \epsilon)$   
 $\nabla_{\theta_\pi} \mathbb{E}_{a \sim \pi_{\theta_\pi}} Q(x, a; \theta_Q) = \mathbb{E}_\epsilon \nabla_x Q(x, \psi(x; \theta_\pi, \epsilon); \theta_Q)$

**EntropyReg:**  $J_\lambda(\theta) = J(\theta) + \lambda H(\pi_\theta) \uparrow$  explor

**SAC:** variant of DDPG/TD3 for  $H$  reg MDPs

**MB DipRL (approx dynamcs model  $f \approx p, r$ )**  
 Init  $\pi$ , data (or  $\{\}$ ), for epis: 1) use  $\pi$ , get data 2) lern  $f, r$  from data, 3) plan new  $\pi$  on estim

**Plang:** cont full obs  $X$ , nonlin trans, constr

**DetDyn:**  $x_{t+1} = f(x_t, a_t)$ , finit  $H$ , at  $t$  mxze:  $J_H(a_{t:t+H-1}) = \sum_{\tau=t:t+H-1} \gamma^{\tau-t} r(x_\tau(a_{t:\tau-1}), a_\tau)$   
 $x_\tau(a_{t:\tau-1}) = f(f(\dots(f(x_t, a_t), a_{t+1}) \dots))$  do  $a_t$ , replan. Opt via  $\nabla$  meth for diff  $r, f$ , cont  $A$  (local min, vanish/explod  $\nabla$ )  $\rightarrow$  use rand shoot. *Rand shot:* Gen rand  $a_{(i), t:t+H-1}$ , pick  $i^* = \text{amax}_i J_H(a_{(i), t:t+H-1})$ . If fin H, sparse  $r$ :  $MPC$  *ValEstim:*  $J_H(a_{t:t+H-1}) = \sum_{\tau=t:t+H-1} \gamma^{\tau-t} r_\tau(x_\tau(a_{t:\tau-1}), a_\tau) + \gamma^H V(x_{t+H})$ , **StocDyn:**  $\max_{a_{t:t+H-1}} \mathbb{E}_{x_{t+1:t+H-1} \sim f(\cdot, a_{t:t+H-1})} [\sum_{\tau=t:t+H-1} \gamma^{\tau-t} r_\tau + \gamma^H V(x_{t+H})]$

$\mathbb{E}$  via *MC traj smpling*, unbias estm of  $J_H$ , aprox via smpl avg. **Param  $\pi$ :** ( $H = 0 \rightarrow$  DDPG)  
 $J_H(\theta) = \mathbb{E}_{x_0 \sim \mu} [\sum_{\tau=0:H-1} \gamma^\tau r_\tau + \gamma^H Q(x_H, \pi(x_H, \theta))] | \theta]$

**UnknownDyn:** follow  $\pi$ , learn  $f, r, Q$  off- $\pi$  from replay buf, replan  $\pi$  based on  $f, r, Q$ . Point est poor perf, err compound  $\rightarrow$  use *BayL*: Model distrib over  $f$  (BNN, GP), use aprox inference (exact, VI, MC..) **Greedy exploi:** 1)  $D = \{\}$ , prior  $P(f | \{\})$  2) repeat: plan new  $\pi$  to  $\max_\pi \mathbb{E}_{f \sim P(\cdot | D)} J(\pi, f)$ , use  $\pi$ , add new data to  $D$ , update post  $P(f | D)$  **PETS algo:** ensmbl of NNs to pred cont Gaus trans distr, MPC for plang. **Explor:** add noise/\* **Thom Smpling:** Like **Greedy** but in 2) smpl model  $f \sim P(\cdot | D)$  then  $\max_\pi J(\pi, f)$  Use epist noise,  $\uparrow$  explor. **Opt explor:** Like **Greedy** but in 2)  $\max_\pi \max_{f \in M(D)} J(\pi, f)$ ; w\  $M(D)$  set of plausible models given  $D$ .