

Summary - Probabilistic Artificial Intelligence

Interpretation of probabilities: Frequentist ($P(A)$) is relative frequency of A in repeated experiments; Bayesian ($P(A)$) is "degree of belief" that A will occur).

$$\text{Bayes rule: } P(X|Y) = \frac{P(X) P(Y|X)}{\sum_{x=x} P(X=x) \cdot P(Y|X=x)} \rightarrow P(\theta|D) = \frac{P(D|\theta) p(\theta)}{\int p(D|\theta) p(\theta) d\theta}$$

Variance: $\text{Var}(X) = E[(X - E[X])^2] = E[X^2] - E[X]^2$ → Bayes rule via Likelihood and prior

$$\text{Var}(aX+bY+c) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X,Y)$$

$$\text{Var}(aX-bY+c) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) - 2ab \text{Cov}(X,Y)$$

$$\text{Cov}(X,Y) = E[(X - E[X]) \cdot (Y - E[Y])].$$

PDF after change of variable: Given r.v. $Y = g(X)$, $g: \mathbb{R} \rightarrow \mathbb{R}$ diff.able injection.

Then $f_Y(y) = f_X(g^{-1}(y)) \cdot \left| \frac{d}{dx}(g^{-1}(y)) \right| = f_X(g^{-1}(y)) \cdot |\det D(g^{-1}(y))|$ ↳ Jacobian in higher dimensions.

Law of unconscious statistician: X r.v., $g: \mathbb{R}^d \rightarrow \mathbb{R}$ diff.able, bijective.

Then $E[g(X)] = \int_{\mathbb{R}^d} g(x) f_X(x) dx = \sum_x g(x) \cdot p(x).$

Law of total expectation: $E[E[X|Y]] = E[X].$

Law of total variance: $\text{Var}(X) = E[\text{Var}(X|Y)] + \text{Var}(E[X|Y])$

Multivariate Gaussian distribution: Let $X \sim \mathcal{N}(\mu, \Sigma)$, then ($X \in \mathbb{R}^d$)

$$p(x) = \frac{1}{(2\pi)^{n/2} \sqrt{|\Sigma|}} \cdot \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

⇒ X, Y Gaussian are independent ⇔ X, Y are uncorrelated.

Let $A = \{i_1, \dots, i_k\} \subseteq \{1, \dots, d\}$, $X = (X_1, \dots, X_d)$ r.v. For the marginal distribution we write $X_A = (X_{i_1}, \dots, X_{i_k}) \sim \mathcal{N}(\mu_A, \Sigma_{AA}).$

⇒ Let $X \sim \mathcal{N}(\mu_X, \Sigma_X)$, $Y \sim \mathcal{N}(\mu_Y, \Sigma_Y)$, M a matrix. Then:

$$X+Y \sim \mathcal{N}(\mu_X+\mu_Y, \Sigma_X+\Sigma_Y) \text{ and } MX \sim \mathcal{N}(M \cdot \mu_X, M \Sigma_X M^T)$$

(→ for X, Y independent.)

Conditional Gaussian: $X = (X_1, \dots, X_d) \sim \mathcal{N}(\mu, \Sigma)$ r.v. (Gaussian random vector), $A, B \subseteq \{1, \dots, d\}$ disjoint. Then define the conditional distr.

$$X_A | (X_B = x_B) \sim \mathcal{N}(\mu_{A|B}, \Sigma_{A|B})$$

where $\mu_{A|B} = \mu_A + \Sigma_{AB} \Sigma_{BB}^{-1} (x_B - \mu_B)$, $\Sigma_{A|B} = \Sigma_{AA} - \Sigma_{AB} \Sigma_{BB}^{-1} \Sigma_{BA}$. → independent of conditional value!!

→ For X, Y jointly Gaussian, X can be viewed as linear function of Y with added random noise. (look at distr. of $X|Y=y$).

Linear regression: Interpolate $y = X \cdot w$, $X \in \mathbb{R}^{n \times d}$, $y, w \in \mathbb{R}^n$, w : weights. → full rank

$$\text{Least squares} = \underset{w}{\operatorname{argmin}} \|y - Xw\|_2^2 = (X^T X)^{-1} X^T y$$

$$\text{Ridge}(\lambda) = \underset{w}{\operatorname{argmin}} \|y - Xw\|_2^2 = (X^T X + \lambda I)^{-1} X^T y$$

Maximum Likelihood Estimation: For observations $(x_1, y_1), \dots, (x_n, y_n)$ we regard the model $y_i = w^T x_i + \varepsilon_i$, where $\varepsilon_i \sim N(0, \sigma_n^2) \rightarrow$ homoscedastic errors. We can equivalently write: $y_i | x_i, w \sim N(w^T x_i, \sigma_n^2)$. Based on this likelihood, we can compute the MLE:
 $\hat{w}_{MLE} = \operatorname{argmax}_w \sum_i \log(p(y_i | x_i, w))$ ($= \hat{w}$ least squares)

Bayesian Inference: Given the model from above, we will now assume that a (Gaussian) prior distribution of our weights $w \sim N(0, \sigma_p^2 I)$. Looking at the conditional version of our problem, the assumption about w allows us to use Bayes-rule to compute the posterior distribution of the weights: $w | x_{1:n}, y_{1:n} \sim N(\mu, \Sigma)$ where $\mu = \sigma_n^{-2} \sum x_i^T y$, $\Sigma = (\sigma_n^{-2} X^T X + \sigma_p^{-2} I)^{-1}$
 \rightarrow Using MLE could give us an estimator for w .

Maximum a Posteriori Estimation: Maximizing the a posteriori probability of our weights (i.e. $p(w | x_{1:n}, y_{1:n})$), results in the MAP estimator: $\hat{w}_{MAP} = \operatorname{argmax}_w \log(p(y_{1:n} | x_{1:n}, w)) + \log(p(w))$

(using Bayes: $p(w | x_{1:n}, y_{1:n}) = p(y_{1:n} | x_{1:n}, w) \cdot p(w) \cdot \text{const.}$, then taking the log)

Notice: $\hat{w}_{MAP} = \operatorname{argmin}_w \|y - Xw\|_2^2 + \frac{\sigma_n^2}{\sigma_p^2} \|w\|_2^2 = \hat{w}_{\text{ridge}} \left(\frac{\sigma_n^2}{\sigma_p^2} \right)$ (last " = ") true
In Gaussian case only!!

To make predictions at a new point x^* , define $f^* = \hat{w}_{MAP}^T x^*$. Then our prediction y^* follows: $y^* | x^*, x_{1:n}, y_{1:n} \sim N(f^*, \sigma_n^2)$

\rightarrow But note that \hat{w}_{MAP} gives no information about the uncertainty of its values. (Problematic when unsure about the model)

Bayesian Linear Regression (BLR): Instead of selecting a single estimation w , we make use of the full posterior distribution of w . For example let $w \sim N(0, I)$ (prior), and the likelihood is $y | x, w, \sigma_n \sim N(w^T x, \sigma_n^2)$. Then the posterior is given by:

$$w | X, y \sim N(\bar{\mu}, \bar{\Sigma}) \quad \text{where } \bar{\mu} = (X^T X + \sigma_n^{-2} I)^{-1} X^T y \quad \text{and } \bar{\Sigma} = (X^T X + \sigma_n^{-2} I)^{-1}.$$

To make a prediction at x^* , define $f^* = w^T x^*$, then $f^* | x^*, x_{1:n}, y_{1:n} \sim N(\bar{\mu}^T x^*, (\bar{\Sigma})^{-1} x^*)$.

Accounting for the noise in the labels, we get the prediction $y^* | x^*, x_{1:n}, y_{1:n} \sim N(\bar{\mu}^T x^*, (\bar{\Sigma})^{-1} x^* + \sigma_n^2)$

We can write $p(y^* | x^*, x_{1:n}, y_{1:n}) = \int p(y^* | x^*, w) p(w | x_{1:n}, y_{1:n}) dw$. Thus BLR averages all w according to the posterior distribution. Note: Ridge regression can be viewed as approximating the posterior by putting all its mass on \hat{w} .

Uncertainty: The uncertainty of a prediction y^* (i.e. its variance) is split into two parts:

$$y^* | X, y, x^* \sim N(\bar{\mu}^T x^*, \underbrace{(\bar{\Sigma})^{-1} x^*}_{\text{epistemic}} + \underbrace{\sigma_n^2}_{\text{aleatoric}})$$

epistemic: uncertainty about the model, due to lack of data. (Uncertainty about f^*)
aleatoric: uncertainty due to irreducible noise

\hookrightarrow uncertainty about y^* given f^* .

BLR hyperparameters: In BLR, we need to specify the variance of the prior σ_p^2 and of the noise σ_n^2 . One simple method to choose them is: (1) Choose $\hat{\lambda} = \sigma_n^2 / \sigma_p^2$ from WMAP via cross-validation, (2) estimate $\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{w}^\top x_i)^2$ and (3) solve for $\hat{\sigma}_p^2 = \hat{\sigma}_n^2 / \hat{\lambda}$. → alternatively: use marginal likelihood (later) to find σ_p^2, σ_n^2

Graphical Models: Represent dependencies of random variables x_1, \dots, x_n with a directed acyclic graph in the following way: $P(x_{1:n}) = P(x_1) \cdot P(x_2 | x_1) \cdot \dots \cdot P(x_n | x_{1:n-1}) = \prod_{i=1}^n P(x_i | x_{S_i})$ → not all prev. are dependent. example: $P(A, B, C, D, E) = P(A) \cdot P(B) \cdot P(C|A, B) \cdot P(D|C) \cdot P(E|C)$ $S_i \subseteq \{1, \dots, i-1\}$

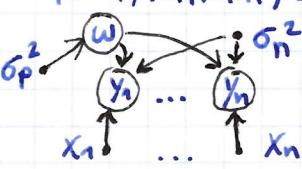


For the Graphical Model of BLR, remember:

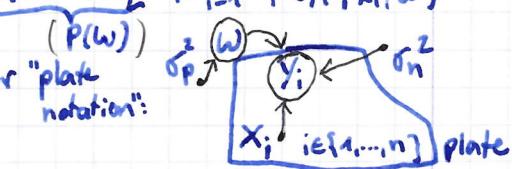
$$\text{prior } p(w) = \mathcal{N}(0, \sigma_p^2 I), P(y_{1:n} | x_{1:n}, w) = \prod_{i=1}^n P(y_i | w, x_i).$$

The joint distribution $p(w, y_{1:n} | x_{1:n}) = p(w | x_{1:n}) \cdot \prod_{i=1}^n P(y_i | x_i, w)$

defines the graph:



($P(w)$)
or "plate notation":



Recursive Bayesian updates: Given new data coming in, we can make our old posterior distribution of the weight, to be ($=P(y_i | x_i, w)$) the new prior (w.r.t. the new data).

Assume $p(w)$ is our prior, and $y_{1:n}$ observations s.t. $P(y_{1:n} | w) = \prod_{i=1}^n P(y_i | w)$

Then we define $p^{(j)}(w) := P(w | y_{1:j})$ → posterior of first j observations, $p^{(0)}(w) = p(w)$.

and recursively $p^{(j)}(w) = p(w | y_{1:j}) = \frac{1}{2} p(w) P_1(y_1 | w) \cdot \dots \cdot P_{j-1}(y_{j-1} | w) \cdot P_j(y_j | w)$
 $= f(p(w | y_{1:j-1}), y_j)$ (for some function f) $= p^{(j-1)}(w)$

→ Advantage: easy to compute; do not need to store all the data at once. ↗ just update model iteratively.

Kernel Trick: Like in BLR, consider the model $f = X w$. Instead of considering a prior over the weights (like in BLR), we can equivalently impose a prior directly on the values of our model f , given our observations X ("Function space view").

Thus: $f | X \sim \mathcal{N}(X E[w], X \text{Var}(w) X^\top) := \mathcal{N}(0, \sigma_p^2 X X^\top)$

↗ note: Cov-matrix is now $n \times n$ whereas in BLR it is $d \times d$ (#param).
 ↗ storage cannot explode from too many parameters.

But notice that data points only enter as inner products.
 This allows us the kernelize f , i.e. write $f | X \sim \mathcal{N}(0, K_{X,X})$

where $K_{X,X} = \begin{pmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{pmatrix}$ with $k(x, x') = \phi(x)^\top \phi(x')$.

We call k a kernel function, and this method the kernel trick.
 → The choice of kernel implicitly determines the class of functions that f is sampled from.

"function space view"

Gaussian Processes (GP): Extending the above idea to infinite domains gives us GP's: Let (x_i) be (infinitely) many random variables,

$X = \{1, \dots, n, \dots\}$ their indices, $M: X \rightarrow \mathbb{R}$ and $k: X \times X \rightarrow \mathbb{R}$ functions.

We write $f \sim GP(M, k)$, if $\forall x \in X$ a random variable f_x s.t.

$\forall A = \{x_1, \dots, x_m\} \subseteq X$ it holds $f_A = [f_{x_1}, \dots, f_{x_m}] \sim \mathcal{N}(M_A, K_{AA})$

↗ jointly normal

where $K_{AA} = \begin{pmatrix} k(x_1, x_1) & \dots & k(x_1, x_m) \\ \vdots & \ddots & \vdots \\ k(x_m, x_1) & \dots & k(x_m, x_m) \end{pmatrix}$

and $M_A = \begin{pmatrix} M(x_1) \\ \vdots \\ M(x_m) \end{pmatrix}$, M : mean function.

↓ k : covariance (kernel) function

Like previously, the choices for k (and μ) implicitly parameterize f . In applications we are typically interested in marginals, i.e.

$$f(x) \sim \mathcal{N}(\mu(x), k(x, x)) \quad (\text{for a specific } x).$$

Covariance (kernel) function: For any function f , and $x, x' \in \mathbb{R}$, we define the covariance function to be

$$k(x, x') = \text{Cov}(f(x), f(x'))$$

if: (1) (Symmetry) $k(x, x') = k(x', x) \forall x, x'$, and
 (2) K_{AA} is positive-semi-definite for any $A \subseteq X$. } see notation from def. of GP.

These conditions ensure that K_{AA} is a valid covariance matrix.
 Intuition: $k(x, x')$ describes how $f(x)$ and $f(x')$ are related.

Composition rules: Let k_1, k_2 be covariance functions on X .

Then the following are valid covariance functions

$$(1) \quad k(x, x') = k_1(x, x') + k_2(x, x')$$

$$(2) \quad k(x, x') = k_1(x, x') \cdot k_2(x, x')$$

$$(3) \quad k(x, x') = C \cdot k_1(x, x') \quad \text{for } C > 0.$$

$$(4) \quad k(x, x') = f(k_1(x, x')) \quad \text{for } f = \exp(\cdot) \text{, or } f = \text{polynomial with pos. coef's.}$$

We call $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ stationary, if $k(x, x') = \tilde{k}(x - x')$ (for some fkt \tilde{k})
 (aka. shift-invariant)

We call $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ isotropic, if $k(x, x') = \tilde{k}(\|x - x'\|_2)$ (for some \tilde{k}).
 Clearly isotropic \Rightarrow stationary.

Examples of Covariance functions:

• Linear kernel: $k(x, x') = x^T x'$. (\rightarrow GP with linear kernel = BLR)

• RBF ((squared exponential/Gaussian) kernel: $k(x, x') = \exp(-\|x - x'\|_2^2/h^2)$
 h: "bandwidth", h small \rightarrow f more complex, h large \rightarrow f simple (less smooth)
 (aka. length scale)

• Exponential kernel: $k(x, x') = \exp(-\|x - x'\|_2/h)$
 h: "bandwidth", h small \rightarrow f complex, h large \rightarrow f simple (generally non-smooth)

• Matérn kernel: $k(x, x') = \frac{2^{1-\gamma}}{\Gamma(\gamma)} \left(\frac{\sqrt{2}\|x - x'\|_2}{\rho} \right)^\gamma \cdot K_\gamma \left(\frac{\sqrt{2}\|x - x'\|_2}{\rho} \right)$
 With parameter γ , Gamma function Γ , modified Bessel fkt K_γ
 and bandwidth (length scale) parameter ρ .
 (continuous, but not diff. able)

Note: $\gamma=1/2 \rightarrow$ exponential kernel ; ($\gamma \rightarrow \infty$) \rightarrow Gaussian (RBF) kernel.

Predictions with GPs: Let $f \sim \text{GP}(\mu, k)$, and $A = \{x_1, \dots, x_m\}$.

Given observations $y_i = f(x_i) + \varepsilon_i$ with $\varepsilon_i \sim \mathcal{N}(0, \sigma_n^2)$.

Then our prediction at point x is distributed as follows:

$$f(x) | x_1, \dots, x_m, y_1, \dots, y_m \sim \text{GP}(\mu'(x), k'(x)) \quad \begin{matrix} \rightarrow \text{etwas unsauber} \\ \rightarrow \text{posterior.} \end{matrix}$$

$$\text{where } \mu'(x) = \mu(x) + k_{x,A} \cdot (K_{AA} + \sigma_n^2 I)^{-1} \cdot (y_A - \mu_A)$$

$$k'(x, x') = k(x, x') - k_{x,A} (K_{AA} + \sigma_n^2 I)^{-1} (k_{x',A})^T \quad \begin{matrix} \rightarrow \text{posterior covariance} \\ \rightarrow \text{doesn't depend on } y \end{matrix}$$

$$\text{with } k_{x,A} = [k(x, x_1), \dots, k(x, x_m)]$$

\rightarrow closed form, i.e. we have again a GP.

Sample from a GP using Cholesky-decomposition: To generate samples from $x \sim \mathcal{N}(\mu, K)$, compute L from $K = LL^T$ (using lin. sy. of eq.), then generate $\varepsilon \sim \mathcal{N}(0, I)$, then compute $x = \mu + L \cdot \varepsilon$.
 (since: $E[x] = \mu + L \cdot E[\varepsilon] = \mu$, $\text{Var}(x) = L \text{Var}(\varepsilon) L^T = LL^T = K$.)

Forward Sampling: Obtain posterior sample-by-sample, i.e.

$$\text{i.e. } f_1 \sim p(f_1), f_2 = p(f_2 | f_1), f_3 = p(f_3 | f_1, f_2), \dots$$

Kalman filter as a GP: First, a Kalman Filter is a dynamical system describing a moving particle: $X_{t+1} = X_t + \varepsilon_t$, $\varepsilon_t \sim \mathcal{N}(0, \sigma_x^2)$. But we only observe perturbed measurements $y_t = X_t + n_t$, $n_t \sim \mathcal{N}(0, \sigma_n^2)$. Then $X_t | y_{1:t} \sim \mathcal{N}(\mu_t, \sigma_t^2)$ and $X_{t+1} | y_{1:t} \sim \mathcal{N}(\mu_{t+1}, \sigma_{t+1}^2 + \sigma_x^2)$. For the Kalman gain k_{t+1} it holds that $X_{t+1} | y_{1:t+1} \sim \mathcal{N}(\mu_{t+1}, \sigma_{t+1}^2)$ where $\mu_{t+1} = \mu_t + k_{t+1} (y_{t+1} - \mu_t)$ and $\sigma_{t+1}^2 = (1 - k_{t+1}) (\sigma_x^2 + \sigma_n^2)$. The Kalman gain k_t is a measure of how much information we gain from a new observation.

A Kalman filter can be seen as a GP: $f(t) = X_t$, $X_0 \sim \mathcal{N}(0, \sigma_0^2)$, $X_{t+1} = X_t + \varepsilon_t$, $\varepsilon_t \sim \mathcal{N}(0, \sigma_x^2)$, then $f \sim \text{GP}(0, K_{KF})$ where $K_{KF}(t, t') = \sigma_0^2 + \sigma_x^2 \min(t, t')$. (\rightarrow Wiener Process kernel)

Optimizing Kernel hyperparameters: A common approach is to choose the parameters Θ , such that the performance of our model is maximized on a validation set. But this approach collapses the uncertainty in f into a point estimate.

↓ Maximize marginal likelihood: (of the training data), is the Bayesian way of doing parameter selection. This way we optimize the effects of Θ across all realizations of f , namely

$$\hat{\Theta}_{MLE} = \operatorname{argmax}_{\Theta} P(Y_{1:n} | X_{1:n}, \Theta) = \operatorname{argmax}_{\Theta} \int P(Y_{1:n}, f | X_{1:n}, \Theta) df \quad \begin{matrix} \rightarrow \text{we get this by} \\ \text{conditioning} \\ \text{on } f. \end{matrix}$$

$$= \operatorname{argmax}_{\Theta} \underbrace{\int P(Y_{1:n} | X_{1:n}, f, \Theta)}_{(1) \text{ likelihood}} \cdot \underbrace{P(f | \Theta)}_{(2) \text{ prior}} df$$

This method typically avoids overfitting (too complex Θ), since this means typically that (1) is small for "most" f and large only for "few" f , while (2) is small. This does not tend to maximize the above integral. The same holds for underfitting (Θ too simple), where (1) is small for "almost all" f , while (2) is large.

The "just right" complexity of Θ is the case when (1) is moderately high for "many" f , and (2) also is moderate. \rightarrow typically maximizes integral.

Remember that $Y_{1:n} | X_{1:n}, \Theta \sim \mathcal{N}(0, K_f(\Theta) + \sigma_n^2 I)$. Define: $K_y(\Theta) = K_f(\Theta) + \sigma_n^2 I$. Then $\hat{\Theta}_{MLE} = \operatorname{argmax}_{\Theta} -\log(P(Y_{1:n} | X_{1:n}, \Theta)) = \underbrace{\frac{1}{2} y^T K_y(\Theta)^{-1} y}_{(1)} + \underbrace{\frac{1}{2} \log(\det(K_y(\Theta)))}_{(2)} + \underbrace{\frac{1}{2} \log(2\pi)}_{\text{const}}$

(1): "goodness of fit", (2): complexity of f 's.

\rightarrow this encourages the trading between large likelihood and large prior.

The gradient of the objective function can be expressed in closed form:

$$\frac{\partial}{\partial \Theta_j} \log(P(Y | X, \Theta)) = \frac{1}{2} \cdot \operatorname{tr}\left((\alpha \alpha^T - K^{-1}) \frac{\partial K}{\partial \Theta_j}\right), \text{ where } \alpha = K^{-1} y, K = K_y(\Theta).$$

\rightarrow this optimization problem is in general non-convex.

\Rightarrow Maximizing marginal likelihood is example of an empirical Bayes method (i.e. estimating a prior distribution from data).

Alternatively, to select a model, we could additionally place a (hyper-) prior on Θ . This hyperprior could then be used to compute the MAP-estimate $\hat{\Theta}_{MAP}$.

Computational issues: To learn a GP, we need to invert matrices $\rightarrow O(n^3)$.

Compared to BLR, which is $O(nd^2)$ (d : Feature dimension).

Thus, it is worthwhile to approximate GP.

Such methods are: Exploiting parallelism, Local GP methods, Kernel function approximations, Inducing point methods.

Exploiting parallelism: The solving of linear systems in GP inference could be sped-up by using certain algorithms on multi-core (GPU) hardware. But note that this doesn't address the cubic scaling in n .

Local methods: Approximate GP by conditioning only on close observations. Notice that as points x, x' are further away, their covariance $k(x, x')$ decays. Thus, when making a prediction at point x , we only condition on the points x' , where $|k(x, x')| \geq \epsilon$ ($\epsilon > 0$). But this method is still computationally quite expensive if "many" points are close by.

Kernel function approximation: Approximate GPs by approximating the kernel function. The idea is to construct a "low-dim." feature map $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^m$, s.t. $k(x, x') \approx \phi(x)^T \phi(x')$.

Then applying BLR has a comp. cost $\Theta(n m^2 + m^3)$.

A such method is Random Fourier Features. But first note that a stationary kernel (i.e. $k(x, x') = k(x - x')$) has a Fourier transform:

$$k(x - x') = \int_{\mathbb{R}^d} p(w) \cdot \exp(j w^T (x - x')) dw \quad j: \text{complex unit.}$$

Theorem (Bochner): k stationary kernel is $\Leftrightarrow p(w) \geq 0$ (nonnegative)
pos. semi definite

→ Thus we can scale the data s.t. $p(w)$ is a probability distribution.
In this case we call $p(w)$ the spectral density of the kernel k .
In the case of the Thm.

For stationary kernels $k(x - x') = k(t)$, the Fourier transform is given by $p(w) = \frac{1}{2\pi} \int_{\mathbb{R}^d} k(t) \cdot \exp(-j \cdot w^T t) dt$

Some examples:

- Gaussian $k(t) = \exp(-\|t\|_2^2/2)$, $p(w) = \frac{1}{(2\pi)^d/2} \cdot \exp(-\|w\|_2^2/2)$
- Exponential: $k(t) = \exp(-\|t\|_1)$, $p(w) = \prod_{j=1}^d 1/\pi(1+w_j^2)$
- Cauchy: $k(t) = \prod_{j=1}^d 2/(1+t_j^2)$, $p(w) = \exp(-\|w\|_1)$.

Key idea of Random Fourier Features (RFF): interpret kernel k as an expectation:

$$k(x - x') = \int_{\mathbb{R}^d} p(w) \cdot \exp(j w^T (x - x')) dw = E_{w \sim p, b \sim \text{Unif}[0, 2\pi]} [Z_{w,b}(x) \cdot Z_{w,b}(x')]$$

$$\text{where } Z_{w,b}(x) = \sqrt{2} \cdot \cos(w^T x + b) \approx \frac{1}{m} \sum_{i=1}^m Z_{w(i), b(i)}(x) \cdot Z_{w(i), b(i)}(x')$$

Random Fourier Features (RFF): Let $k(x, x') = k(x - x')$ be a pos. def. stationary kernel. RFF constructs $Z: \mathbb{R}^d \rightarrow \mathbb{R}^m$ s.t. $k(x, x') \approx Z(x)^T Z(x')$.

(1) Construct Fourier Transform of the kernel:

$$P(w) = \frac{1}{2\pi} \int_{\mathbb{R}^d} k(t) \cdot \exp(-j w^T t) dt$$

(2) Randomly sample $w_1, \dots, w_m \sim P$, and $b_1, \dots, b_m \sim \text{Uniform}([0, 2\pi])$

(3) Define: $Z(x) := \sqrt{2/m} \cdot [\cos(w_1^T x + b_1) + \dots + \cos(w_m^T x + b_m)]$

Theorem: $M \subseteq \mathbb{R}^d$ compact. diam(M) its diameter. $Z(x)$ as in RFF:

$$P(\sup_{x, x' \in M} |Z(x)^T Z(x') - k(x, x')| \geq \epsilon) \leq 2^d \left(\frac{\text{op} \cdot \text{diam}(M)}{\epsilon} \right)^2 \cdot \exp\left(-\frac{m \epsilon^2}{4(d+2)}\right) \quad (\text{op}^2 = E[w^T w]).$$

→ BLR with feature map Z approximates a GP with stationary kernel.
→ Error probability decays exponentially fast. → in the dimension of the Fourier Feature Space.

Problem: Fourier Features can be wasteful, since they approximate the kernel function uniformly well (across all x). But we only need accurate approx. for the points in training/test data.

Inducing points method: Approximate GP's by only considering a (random) subset of the training data during learning.

One subsampling method is the inducing points method. The idea is to summarize the training data around so-called inducing points.

Write $U = \{\bar{x}_1, \dots, \bar{x}_k\}$ (treat inducing points as hyperparam.).

Let $f = [f(x_1), \dots, f(x_n)]^T$ be the evaluations of our training data \bar{X} ,

$f^* = f(\bar{x}^*)$ the prediction at some evaluation point $\bar{x}^* \in \mathcal{X}$.

Further denote with $u = [f(\bar{x}_1), \dots, f(\bar{x}_k)]^T \in \mathbb{R}^k$ the predictions of our model f at the inducing points U .

Then our joint prior is given by

$$p(f, f^*) = \int_{\mathbb{R}^k} p(f, f^*, u) du = \int_{\mathbb{R}^k} p(f, f^* | u) \cdot p(u) du$$

$$\approx \int_{\mathbb{R}^k} p(f^* | u) \cdot p(f | u) \cdot p(u) du.$$

→ by using marginalization

The key idea is to approximate the joint prior, by assuming that f and f^* are conditionally independent given u . We call "Training conditional": $p(f | u) = \mathcal{N}(K_{f,u} \cdot K_{u,u}^{-1} \cdot u, K_{ff} - Q_{f,f})$

"Testing conditional": $p(f^* | u) = \mathcal{N}(K_{f^*,u} \cdot K_{u,u}^{-1} \cdot u, K_{f^*,f^*} - Q_{f^*,f^*})$

$$\text{where } Q_{ab} = K_{a,u} \cdot K_{u,u}^{-1} \cdot K_{u,b}.$$

→ intuitively: K_{AA} (prior cov.), Q_{AA} (cov. "explained" by inducing points)

→ But computing this covariance matrix is expensive. Two approximations to the covariance of the training/testing conditional are Subset of Regressors (SoR), and Fully indep. training cond. (FITC).

Subset of Regressors (SoR): SoR approximation simply leaves out the covariance of the training/testing conditionals:

$$q_{\text{SoR}}(f | u) = \mathcal{N}(K_{f,u} \cdot K_{u,u}^{-1} \cdot u, 0) \text{ and } q_{\text{SoR}}(f^* | u) = \mathcal{N}(K_{f^*,u} \cdot K_{u,u}^{-1} \cdot u, 0).$$

→ resulting model is degenerate GP with $k_{\text{SoR}}(x_i, x') = k(x_i, u) \cdot K_{u,u}^{-1} \cdot k(u, x')$.

→ computing costs dominated by inverting $K_{u,u}^{-1}$. Thus cubic in K (#inducing points), but linear in n (#data points).

Fully independent training conditional (FITC): FITC approximation only keeps variances, and forgets about covariance:

$$q_{\text{FITC}}(f | u) = \prod_{i=1}^n p(f_i | u) = \mathcal{N}(K_{f,u} \cdot K_{u,u}^{-1} \cdot u, \text{diag}(K_{f,f} - Q_{f,f}))$$

$$q_{\text{FITC}}(f^* | u) = \mathcal{N}(K_{f^*,u} \cdot K_{u,u}^{-1} \cdot u, \text{diag}(K_{f^*,f^*} - Q_{f^*,f^*}))$$

→ computing costs dominated by inverting $K_{u,u}^{-1}$. Thus cubic in k (#inducing points), but only linear in n (#data points)

How to pick inducing points: Choose randomly, greedily (according to e.g. variance), (random) equally spaced, or treat u as hyperparameters, and maximize marginal likelihood of the data.

→ It is most important to ensure that u is representative of the test and training (!) data.

Bayesian Learning: More generally BL works the following way:

Assume a prior $p(\theta)$ (e.g. $\mathcal{N}(\theta; \sigma_p^2 I)$). The

likelihood is: $P(Y_{1:n}|X_{1:n}, \theta) = \prod_{i=1}^n p(y_i|X_i, \theta)$. $\rightarrow \ln \text{BLR} \sim \mathcal{N}(\theta^T X, \sigma_n^2)$

Then posterior: $p(\theta|X_{1:n}, Y_{1:n}) = \frac{1}{Z} p(\theta) \prod_{i=1}^n p(y_i|X_i, \theta)$

$$\text{where } Z = \int p(\theta) \prod_{i=1}^n p(y_i|X_i, \theta) d\theta$$

Predictions are given by: $p(y^*|x^*, X_{1:n}, Y_{1:n}) = \int p(y^*|x^*, \theta) \cdot p(\theta|X_{1:n}, Y_{1:n}) d\theta$

\rightarrow For BLR and ~~GP regression~~, these integrals are closed-form.

\rightarrow In general the integrals are not closed-form \rightarrow need approximations!

can be
computed
directly.

To approximate intractable distributions $p(\theta|y) = \frac{1}{Z} p(\theta, y)$

we will assume we can evaluate the joint distribution $p(\theta, y)$,

but not the normalizer Z . General approaches focus on approximating $p(\theta|y)$ as a whole (e.g. Variational Inference, Laplace approx.)

Laplace Approximation: We approximate the posterior $p(\theta|X_{1:n}, Y_{1:n}) \approx q(\theta)$

using a Gaussian approximation (i.e. 2nd-order Taylor expansion) of the log-posterior around its mode $\hat{\theta}$. Then the L-approx. is:

$$q(\theta) = \mathcal{N}(\hat{\theta}, \Lambda^{-1}) \quad \text{where } \hat{\theta} = \underset{\theta}{\operatorname{argmax}} P(\theta|y), \quad \Lambda = -\nabla \nabla \log(p(\hat{\theta}|y))$$

\rightarrow But this method can lead to poor (e.g. overconfident) approximations, since it first greedily seeks the mode, and then matches the curvature (arbitrarily bad, compared to the true posterior). \rightarrow approx. of a Gaussian is exact

Example: Bayesian logistic regression. Prior $p(w) = \mathcal{N}(0, \sigma_p^2 I)$, and

likelihood $p(Y_{1:n}|X_{1:n}, w) = \text{Bernoulli}(\sigma(w^T x))$, with $\sigma(z) = 1/(1 + \exp(-z))$ the logistic function. To find the approximate posterior, first compute the mode: $w = \underset{w}{\operatorname{argmax}} p(w|Y_{1:n}) = \underset{w}{\operatorname{argmax}} \frac{1}{Z} p(w) p(Y_{1:n}|w)$

take

$$\log \Rightarrow \hat{w} = \underset{w}{\operatorname{argmin}} \underbrace{\frac{1}{2\sigma_p^2} \|w\|_2^2}_{=: \lambda} + \sum_{i=1}^n \underbrace{\log(1 + \exp(-y_i w^T x_i))}_{\text{"logistic loss": } l(w, x_i, y_i)}$$

"logistic loss": $l(w, x_i, y_i)$

\rightarrow Apply SGD to find \hat{w} , since the above term is simply standard log. regression.

SGD recap: minimize $L(\theta) = \mathbb{E}_{x \sim p}[l(\theta, x)] = \frac{1}{n} \sum_{i=1}^n l(\theta, x_i)$. (1) initialize

θ_1 , (2) draw from mini-batch $B = \{x_1, \dots, x_m\}$, $x_i \sim p$, (3) update

$$\theta_{t+1} \leftarrow \theta_t - \eta_t \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} l(\theta_t, x_i)$$

In our case: (1) initialize w , (2) pick point from (X, Y) uniformly at random

(3) compute prob. of missclassification $\hat{P}(Y=-y|w, x) = 1/(1 + \exp(y w^T x))$

$$(4) \text{ update } w_t \leftarrow w_t (1 - 2\lambda \eta_t) + \eta_t y x \hat{P}(Y=-y|w, x)$$

Second: Compute the variance of our Laplace approximation:

$$\Lambda = -\nabla \nabla \log(p(w|X_{1:n}, Y_{1:n})) = X^T \text{diag}([\pi_i \cdot (1 - \pi_i)]_i) \cdot X, \quad \pi_i = \sigma(\hat{w}^T x_i)$$

$\rightarrow \Lambda$ does not depend on normalization constant Z .

Lastly: For Bayesian Logistic regression, the predictive distribution for a test point x^* is given by

$$p(y^*|x^*, X_{1:n}, Y_{1:n}) \approx \int \delta(y^* f) \cdot W(f; \hat{w}^T x^*, x^{*T} \Lambda^{-1} x^*) df$$

Surprise: Let $u = P(A)$ for an event A . We define the surprise to be $S[u] = -\log(u)$. (note $S[u] \geq 0$).

Axiomatic characterization: (1) $S[1]=0$, (2) $S[u] > S[v] \Rightarrow u < v$,

(3) S is continuous, (4) for independent events: $S[uv] = S[u] + S[v]$.

Entropy: is a notion of uncertainty about the distribution p .

$$H[p] = E_{x \sim p} [-\log(p(x))] = - \int p(x) \log(p(x)) dx = - \sum_x p(x) \log_2(p(x))$$

product distribution: $p(X_1:d) = \prod_{i=1}^d p_i(x_i) \Rightarrow H[p] = \sum_{i=1}^d H[p_i]$

example (Gaussian): $H[\mathcal{N}(\mu, \Sigma)] = d/2 \cdot \log(2\pi e) + \frac{1}{2} \log(\det(\Sigma))$

$\rightarrow H[p]$ larger: more unsure about $X \sim p$.

\hookrightarrow discrete case

d -dimension of gaussian

Kullback-Leibler Divergence: First, define the cross-entropy of a distribution q relative to the distribution p .

$$H[p \parallel q] = E_{x \sim p} [-\log(q(x))] (= H[p] + KL(p \parallel q)) \geq H[p].$$

Intuition: measure of how close q is to the true distribution p . (closer \rightarrow lower). Thus, the KL-divergence is defined to be

$$KL(p \parallel q) = H[p \parallel q] - H[p] = \int p(\theta) \log\left(\frac{p(\theta)}{q(\theta)}\right) d\theta = E_{\theta \sim p} [\log\left(\frac{p(\theta)}{q(\theta)}\right)]$$

\rightarrow this is a measure of distance between the distributions p and q .

Properties:

$$(1) KL(p \parallel q) \geq 0, (2) KL(p \parallel q) = 0 \Leftrightarrow p = q \text{ a.s.}$$

$$(3) \text{In general: } KL(p \parallel q) \neq KL(q \parallel p) \text{ (not symmetric).}$$

\rightarrow KL measures the additional expected surprise when observing samples from p , but assuming the (wrong) distribution q . (q assigns high probability to outcomes which are likely according to p)

\rightarrow minus the surprise inherent in p .

Example: $KL(q \parallel p(\theta)) = \frac{1}{2} \sum_{i=1}^d (\theta_i^2 + \mu_i^2 - 1 - \log(\theta_i^2))$, $q \sim \mathcal{N}(\mu, \text{diag}(\theta_1^2, \dots, \theta_d^2))$, $p \sim \mathcal{N}(0, I)$

\hookrightarrow KL for Bayesian logistic regression

Variational Inference: We approximate the posterior $p(\theta | X_{1:n}, Y_{1:n}) \approx q(\theta)$ by finding a "simple" distribution that approximates p well:

$$q^* \in \operatorname{argmin}_{q \in Q} KL(q \parallel p) = \operatorname{argmin}_{\lambda \in \mathbb{R}^n} KL(q_\lambda \parallel p) \quad \rightarrow \text{quality of approx. is hard to analyze.}$$

Q: "Variational family" (needs to be specified)

\rightarrow "reverse KL-divergence"

Note: Let $q_1^* = \operatorname{argmin}_{q \in Q} KL(p \parallel q)$ and $q_2^* = \operatorname{argmin}_{q \in Q} KL(q \parallel p)$.

When approximating the true posterior p , using q_1^* , q_2^* tends to greedily select the mode, and underestimate the variance. (\rightarrow overconfident) While q_1^* gives a more conservative (usually desired) estimation.

\rightarrow For Q the family of Gaussians: q^* matches 1st and 2nd moments of p .

\rightarrow Minimizing KL-divergence is equivalent to maximum likelihood estimation on an infinitely large sample size.

\rightarrow While for q^* we would prefer forward-KL, variational inference chooses backward-KL for computational reasons (MC: sample from p otherwise). But minimizing reverse-KL still recovers the true posterior a.s..

Evidence Lower Bound (ELBO): Instead of minimizing the reverse-KL in Variational inference, we can maximize the ELBO. (equivalently).

$$\operatorname{argmax}_{q \in Q} KL(q \parallel p(\cdot | X_{1:n}, Y_{1:n})) = \operatorname{argmax}_{q \in Q} E_{\theta \sim q} [\log(p(Y_{1:n} | X_{1:n}, \theta))] - KL(q \parallel p(\cdot))$$

$$= \operatorname{argmax}_{q \in Q} E_{\theta \sim q} [\log(p(Y_{1:n}, \theta | X_{1:n}))] + H[q], \quad \rightarrow \text{prior } p(\theta)$$

$$= L(q, p(\cdot | X, Y))$$

We call $L(q, p(\cdot | X, Y))$ the evidence lower bound (ELBO).

Because: log-evidence: $\log(p(Y_{1:n} | X_{1:n})) \geq ELBO \rightarrow$ use Jensen inequality

\rightarrow ELBO selects q with large joint likelihood $p(Y_{1:n}, \theta | X_{1:n})$ (i.e. data is likely, and q is close to prior), and with large Entropy $H[q]$ (q is simple).

\rightarrow Maximizing ELBO with noninformative prior $p(\theta) \equiv \text{const.}$ is equivalent to MLE.

Gradient of ELBO: To maximize ELBO common methods (e.g. SGD) would require an unbiased estimate of the gradient.

$$\nabla_{\lambda} L(q_{\lambda}, p; X, y) = \nabla_{\lambda} E_{\Theta \sim q_{\lambda}} [\log(p(Y_{1:n} | X_{1:n}, \theta))] - \nabla_{\lambda} KL(q_{\lambda} || p(\cdot))$$

While $\nabla_{\lambda} KL$ can easily be computed for common variational families, $\nabla_{\lambda} E$ is more difficult to compute. Since the expectation integrates over the measure q_{λ} , which depends on the variational parameters λ .
 → Thus we rewrite the gradient such that Monte Carlo sampling becomes possible. (e.g. Reparametrization trick, score gradients).
 → maximizing ELBO using stochastic optimization is only twice as expensive (for diagonal q) as MAP inference. ($\nabla_{\lambda} KL$ is easy).
 But we can achieve better performance: e.g. Natural gradients, Variance red techniques.

Reparameterization Trick: allows a change of "densities" s.t. the expectation above is w.r.t. ϕ , which does not depend on λ .

Let $\Sigma \sim \Phi$ a r.v. (with Φ independent of λ). Given $g: \mathbb{R}^d \rightarrow \mathbb{R}^d$ invertible, diff. able and $\theta = g(\Sigma, \lambda)$. Then it holds that

$$q(\theta | \lambda) = \phi(\Sigma) \cdot |\det(\nabla_{\Sigma} g(\Sigma, \lambda))|^{-1} \text{ and } E_{\Theta \sim q_{\lambda}} [f(\theta)] = E_{\Sigma \sim \Phi} [f(g(\Sigma, \lambda))]$$

→ We obtain stochastic gradients via: $\nabla_{\lambda} E_{\Theta \sim q_{\lambda}} [f(\theta)] = E_{\Sigma \sim \Phi} [\nabla_{\lambda} f(g(\Sigma, \lambda))]$

Example: Bayesian Logistic Regression.

Suppose we use a Gaussian variational approximation $q_{\lambda}(\theta) = \mathcal{N}(\mu, \Sigma)$.
 Let $\Sigma \sim \mathcal{N}(0, I)$, and use reparam. $\theta = g(\Sigma, \lambda) = \Sigma^{1/2} \Sigma + \mu$.

Then $\Sigma = \Sigma^{1/2} (\theta - \mu)$ and $\phi(\Sigma) = q_{\lambda}(\theta) \cdot |\det(\Sigma^{1/2})|^{-1}$.

The ELBO gradient (i.e. $\nabla_{\lambda} E$ -bit) simplifies to: (let $C := \Sigma^{1/2}$)

$$\begin{aligned} \nabla_{\lambda} E_{\Theta \sim q_{\lambda}} [\log(p(Y_{1:n} | X_{1:n}, \theta))] &= \nabla_{C, \mu} E_{\Sigma \sim \mathcal{N}(0, I)} [\log(p(Y_{1:n} | X_{1:n}, \theta))] \Big|_{\theta=C\Sigma+\mu} \\ &\approx n \cdot \frac{1}{m} \sum_{j=1}^m \nabla_{C, \mu} \log(p(Y_j | X_{i,j}, \theta)) \Big|_{\theta=C\Sigma+\mu} \quad (\rightarrow \text{using Monte Carlo Sampling.}) \end{aligned}$$

→ unbiased

Markov Chains: A (stationary, time-homogenous) Markov Chain is a sequence of r.v.'s $(X_n)_{n \in \mathbb{N}}$ with prior density $P(X_1)$, and transition probabilities $P(X_{t+1} | X_t)$ independent of t ($X_{t+1} \perp X_0: t+1 | X_t$: "Markov property")
 As the state space $\{x_1, \dots, x_n\}$ is finite, MC can be described by the transition matrix $P = \begin{pmatrix} P(x_1 | x_1) & \dots & P(x_n | x_1) \\ \vdots & \ddots & \vdots \\ P(x_1 | x_n) & \dots & P(x_n | x_n) \end{pmatrix}$ of transition probabilities.

→ The rows of P sum to 1.

→ P describes a directed graph, which we call transition graph.

We call an MC ergodic, if $\exists t$ finite, s.t. every state can be reached from every state in exactly t steps ($\Rightarrow P^t$ has only strictly positive entries).
 (\Rightarrow irreducible (can reach every state from everywhere) and aperiodic ($\forall k > N: \pi^{(k)}(X_i, X) > 0$))

arbitrary

→ If MC ergodic, then it has a unique stationary distribution

$$\pi(x) = \sum_{x' \in S} P(x|x') \cdot \pi(x') > 0 \quad \text{s.t.} \quad \lim_{N \rightarrow \infty} P(X_N=x) = \pi(x).$$

$$\Rightarrow \pi = \pi \cdot P \quad (\text{for } \pi \text{ a row vector})$$

→ irrespective of $P(x_1)$!

A MC satisfies the detailed balance equation for an unnormalized distribution $Q(x)$, if: $\frac{1}{z} Q(x) \cdot P(x'|x) = \frac{1}{z} Q(x') P(x|x')$

$$Q(x), \text{if: } \frac{1}{z} Q(x) \cdot P(x'|x) = \frac{1}{z} Q(x') P(x|x')$$

normalizing factor.

⇒ The Markov Chain has stationary distribution $\frac{1}{z} Q(x)$.

Markov Chain Monte Carlo (MCMC): We estimate the distribution of the prediction (from Bayesian Learning) directly, by interpreting the integral as an expectation over the posterior:

$$p(y^* | x^*, x_{1:n}, y_{1:n}) = \int p(y^* | x^*, \theta) \cdot p(\theta | x_{1:n}, y_{1:n}) d\theta = E_{\theta \sim p(x_{1:n}, y_{1:n})} [p(y^* | x^*, \theta)]$$

$$\approx \frac{1}{m} \sum_{i=1}^m p(y^* | x^*, \theta^{(i)})$$

→ Error of approximation decreases exponentially in m.

↪ NOT for MCMC!
(since samples are not i.i.d.)

We approximate this term using Monte Carlo sampling of i.i.d. $\theta^{(i)} \sim p(\cdot | x_{1:n}, y_{1:n})$. We sample the $\theta^{(i)}$ using MC.

From the detailed balance equation it follows that we can take an unnormalized distribution and construct an MC which has the normalized distribution as stationary distributions.

→ Note: samples from MC are NOT i.i.d. if from same MC

↪ MC are also very easy to simulate.

Metropolis-Hastings Algorithm: describes a way of constructing a MC with desired stationary distribution $p(x) = \frac{1}{Z} q(x)$. (normalized q).

Given a proposal distribution $r(x'|x)$, which in state x proposes state x' . Then with probability $\alpha(x'|x)$ we follow the proposal, and with probability $1-\alpha$ we stay in state x .

$$\text{Here } \alpha(x'|x) = \min\left\{1, \frac{P(x') \cdot r(x|x')}{p(x) \cdot r(x'|x)}\right\}$$

→ can also replace P with q.

→ the resulting MC has stationary distribution $p(x) = \frac{1}{Z} q(x)$.

→ Many ways to choose proposal distr. r (e.g. Gibbs Sampling).

Gibbs-Sampling: is a Metropolis-Hastings Algorithm where $\alpha(x, x') = 1$, and proposal distribution $r(x'|x)$ s.t. we first choose $i \in \{1, \dots, n\}$ uniformly at random, then set $x_i = [x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n]$, and update x_i by sampling from $P(x_i | x_{-i})$.

→ from components of initial value $x \in \mathbb{R}^n$.

→ This satisfies the detailed balance equation (since choosing i uniformly at random ensures that the MC is ergodic).

→ finds states that are successively more likely.

→ Choosing the i in fixed order still achieves correct stationary distribution, but does not fulfill the detailed balance equation.

Generally sampling $x_i | x_{-i}$ is efficient, since $p(x_i | x_{-i}) = q(x_i | x_{-i}) / \sum_{x_i} q(x_i | x_{-i})$

→ can be efficiently evaluated

Ergodic Theorem: Let $(X_n)_{n \in \mathbb{N}}$ be an ergodic MC over a finite state space S , with stationary distribution π , and $f: S \rightarrow \mathbb{R}$.

$$\text{Then } \frac{1}{n} \sum_{i=1}^n f(X_i) \xrightarrow{n \rightarrow \infty} \sum_{x \in S} \pi(x) \cdot f(x) = E_{x \sim \pi} [f(x)], (x_i \sim X_i | x_{i-1})$$

↪ like strong law of large numbers for MC.

→ thus MCMC works (can approximate expectation), even though the samples from the MC are strongly dependent.

→ Simulating from a single MC yields an unbiased estimator.

→ No information about rate of convergence, and variance.

→ We usually discard the first samples of the MC ("burn-in time").

→ tune how many should be discarded, and sampled in total

MCMC for continuous random variables: Generalizing MCMC for continuous r.v. We focus on positive (so-called: Gibbs) distributions of the form

$$p(x) = \frac{1}{Z} \cdot \exp(-f(x)) \quad f: \text{"energy function".}$$

We call p log-concave, if f is convex. (e.g. Bayesian logistic regression)
The adapted Metropolis-Hastings Algorithm uses

$$\alpha(x' | x) = \min\left\{1, \frac{r(x'|x)}{r(x|x')} \cdot \exp(f(x) - f(x'))\right\} \text{ for } \alpha.$$

Normally distributed proposal: $r(x'|x) = \mathcal{N}(x'; x, \Sigma)$

$$\text{then } \frac{r(x'|x)}{r(x|x')} = 1 \Rightarrow \alpha = \min\{1, \exp(f(x) - f(x'))\}$$

Metropolis adjusted Langevin Algorithm (MALA): describes an MCMC for continuous r.v. with improved proposal distribution:

$$r(x'|x) = N(x'; x - \gamma \nabla f(x), 2\gamma I)$$

or LMC (Langevin Monte Carlo)
locally non-convex is also fine!

If the distr. p (to be approximated) is log-concave (e.g. Bayesian log. regression) then MALA efficiently converges to its stationary distribution.
 → But proposal and acceptance step in MALA require full access to energy function f . → expensive for large datasets.

Stochastic Gradient Langevin Dynamics (SGLD): improves efficiency of MALA using stochastic gradient estimates, decaying step sizes and by skipping the accept/reject step. The algorithm works as follows:

sample $\theta \sim \frac{1}{n} \exp(\log(p(\theta)) + \sum_{i=1}^n \log(p(y_i|x_i, \theta)))$ from Bayesian posterior.

(1) Initialize θ_0 , (2) For $t=0, 1, \dots$ repeat $i_1, \dots, i_m \sim \text{Uniform}\{1, \dots, n\}$,
 $\varepsilon_t \sim N(0, 2\eta_t I)$, update $\theta_{t+1} \leftarrow \theta_t - \eta_t (\nabla_\theta \log(p(\theta_t)) + \frac{n}{m} \sum_{j=1}^m \nabla_\theta \log(p(y_j|\theta_t, x_j)))$

→ Guaranteed convergence to stationary distribution for $\eta_t \in \Theta(t^{-1/3})$. → bounded above and below by $t^{-1/3}$.

choose mini batch.
n: total sample size.

+ ε_t
↓ Gaussian noise.

Using Hamiltonian Monte Carlo (HMC), performance can be improved even more by adding a momentum term to (S)GD.

→ In general we can guarantee convergence to target distribution, but for general distributions convergence/mixing may be slow.

Bayesian Neural Networks: In practice, we often get better performance in Bayesian Learning, when introducing nonlinearity: BNN's.

In BNN's we specify a prior distribution over the weights of the NN $p(\theta)$, and parameterize our likelihood functions via Neural Networks.

Example: Gaussian prior $p(\theta) = N(\theta; 0, \sigma_p^2 I)$, and likelihood $p(y|X, \theta) = N(f(X, \theta), \sigma^2)$ where f is non-linear. → NN approximates f .

→ BNN's learn distribution over weights of the Neural Network.

→ Alternatively BNN's could handle more complex likelihoods, such as $p(y|X, \theta) = N(f_1(x, \theta), \exp(f_2(x, \theta)))$, and BNN would model both mean and (log-) variance as outputs of one NN.

MAP Estimation for BNNs: In the case of heteroscedastic error, we use the likelihood model $p(y|X, \theta) = N(\mu(x, \theta), \sigma(x, \theta))$ μ, σ : Functions.

Then $\hat{\theta}_{\text{MAP}} = \arg \min_{\theta} -\log(p(\theta)) - \sum_{i=1}^n \log(p(y_i|x_i, \theta))$

$$= \arg \min_{\theta} \lambda \|\theta\|_2^2 + \frac{1}{2} \sum_{i=1}^n [\frac{1}{\sigma(x_i, \theta)^2} \|y_i - \mu(x_i, \theta)\|_2^2 + \log(\sigma(x_i, \theta))^2]$$

→ Model can explain y_i through accurate model μ , or variance σ , but is penalized for choosing large variances. Thus reduce losses by attributing some points to large variance. → "learn" aleatoric uncertainty.

→ MAP only makes point estimates → don't model epistemic uncertainty. → lack of data uncertainty.

Approximate Inference for BNNs: Learning and inference in BNNs are generally intractable (even for Gaussian prior and likelihood) when noise is not assumed to be homoscedastic and known.

→ Like before, the integrals from BL are intractable.

→ Can use approximate inference techniques (e.g. Laplace Approx., Black-box stochastic variational inf., SGLD) at similar cost as MAP/SGD.

→ automatic differentiation

→ Computing gradients is not a problem, since well-implemented in NN
 → specialized approximate inference methods for BNNs: Monte-Carlo Dropout, Probabilistic Ensembles.

Now we will describe many such methods for approximate inference.

Variational Inference for BNNs (Bayes by Backprop): As we have done with VI, approximate the posterior $p(\theta | x_{1:n}, y_{1:n}) \approx q(\theta)$ by maximizing the ELBO: $\arg\min_{q \in Q} KL(q || p(\cdot | x_{1:n}, y_{1:n})) = \arg\max_{q \in Q} E_{\theta \sim q} [\log(p(y_{1:n} | x_{1:n}, \theta))] - KL(q || p(\cdot))$

Where we usually choose Q to be the fam. of Gaussians: $q(\theta | \lambda) = \mathcal{N}(\theta; M, \Sigma)$. To maximize, we need to compute the gradient. $\nabla_{\lambda} KL$ can be expressed in closed form (for Gaussians). For $\nabla_{\lambda} E$ we obtain unbiased gradient estimates using the reparametrization trick:

$$E_{\theta \sim q} [\log(p(y_{1:n} | x_{1:n}, \theta))] = E_{\Sigma \sim \mathcal{N}(0, I)} [\log(p(y_{1:n} | x_{1:n}, \theta)) | \theta = \Sigma^{1/2} \varepsilon + \mu]$$

$$\text{where } \Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_d^2), \Sigma^{1/2} = \text{diag}(\sigma_1, \dots, \sigma_d).$$

Then we use backpropagation (automatic differentiation) to compute the gradients of the Likelihood.

$$\text{Inference: } p(y^* | x^*, x_{1:n}, y_{1:n}) = \int p(y^* | x^*, \theta) \cdot p(\theta | x_{1:n}, y_{1:n}) d\theta$$

$$\approx E_{\theta \sim q} [p(y^* | x^*, \theta)] \approx \frac{1}{m} \sum_{i=1}^m p(y^* | x^*, \theta^{(i)}) \rightarrow \text{monte carlo sampling}$$

→ Intuitively, we draw NN weights according to the posterior q_{λ} , and then average-out their predictions.

Using the inference approximation we can get the prediction distribution (for a Gaussian likelihood):

$$p(y^* | x_{1:n}, y_{1:n}, x^*) \approx \frac{1}{m} \sum_{i=1}^m \mathcal{N}(\mu(x^*, \theta^{(i)}), \sigma^2(x^*, \theta^{(i)})).$$

$$\text{Thus approx. mean: } E[y^* | x^*, x_{1:n}, y_{1:n}] \approx \frac{1}{m} \sum_{i=1}^m \mu(x^*, \theta^{(i)}) =: \bar{\mu}(x^*)$$

$$\text{and approx. variance: } \text{Var}(y^* | x_{1:n}, y_{1:n}, x^*) = E[\text{Var}(y^* | x^*, \theta)] + \text{Var}(E[y^* | x^*, \theta])$$

$$\approx \underbrace{\frac{1}{m} \sum_{i=1}^m \sigma^2(x^*, \theta^{(i)})}_{\text{aleatoric uncertainty}} + \underbrace{\frac{1}{m} \sum_{i=1}^m (\mu(x^*, \theta^{(i)}) - \bar{\mu}(x^*))^2}_{\text{epistemic uncertainty}},$$

MCMC for BNNs: Like before, we can use MCMC methods to approximate the prediction directly: $p(y^* | x^*, x_{1:n}, y_{1:n}) = \frac{1}{N} \sum_{i=1}^N p(y^* | x^*, \theta^{(i)})$

Algorithms like SGLD or SG-HMC allow efficient computation of gradients.

→ To avoid "burn-in" period, drop first n samples

→ Typically cannot afford to store all N samples / models. Thus, we need to summarize the iterates. (i.e. combine the $\theta^{(i)}$).

subsampling: only keep track of m snapshots of weights $[\theta^{(1)}, \dots, \theta^{(m)}]$ (acc. to some schedule) and use them for inference.

Gaussian approximation: summarize weights by approx. them with a Gaussian distribution $q(\theta | M_{1:d}, \Sigma_{1:d}^2)$ where

$$M_i = \frac{1}{T} \sum_{j=1}^T \theta_i^{(j)} \quad \text{and} \quad \Sigma_i^2 = \frac{1}{T} \sum_{j=1}^T (\theta_i^{(j)} - M_i)^2$$

→ This can be implemented using running averages.

→ for prediction, simply sample weights from $q(\theta | M_{1:d}, \Sigma_{1:d}^2)$

→ SWAG method: use SGD (without Gaussian noise) to generate $\theta^{(1)}, \dots, \theta^{(T)}$.

weight samples
of one iteration

Dropout as Variational Inference: In dropout regularization, the key idea is to randomly ignore ("drop out") hidden units during training (i.e. during each iteration of SGD) with probability p . Dropout can be viewed as performing variational inference with the variational family: $q(\theta | \lambda) = \prod_j q_j(\theta_j | \lambda_j) = \prod_j (p \delta_0(\theta_j) + (1-p) \delta_{\lambda_j}(\theta_j))$

↓ i.e. each weight θ_j is either set to 0 with prob. p , or set to λ_j with prob. $(1-p)$. 13

$$\delta_a(x) = \begin{cases} 1 & x=a \\ 0 & x \neq a \end{cases}$$

↓ (Dropout as Variational Inference)

→ Thus ordinarily training a NN with dropout is equivalent to performing approximate Bayesian inference.

The approximation of the prediction from Variational inference:

$$P(y^*|x^*, x_{1:n}, y_{1:n}) = \frac{1}{m} \sum_{i=1}^m P(y^*|x^*, \theta^{(i)})$$

Each sample $\theta^{(i)}$ corresponds to a NN with weights given by λ , where each unit is set to 0 with probability p .

→ Thus dropout also performed during prediction! (not only training)

Probabilistic Ensembles: is another approach for approximate inference

for BNNs. We know that VI for BNNs is like averaging the predictions of m NN. Thus, here, we immediately learn the weights of m NNs. For this we choose m training sets randomly, by sampling uniformly from the data, with replacement. Then we obtain m MAP estimates of the weights $\theta^{(i)}$, and approximate the prediction:

$$P(y^*|x^*, x_{1:n}, y_{1:n}) = E_{\theta \sim P(\cdot|x_{1:n}, y_{1:n})} [P(y^*|x^*, \theta)] \approx \frac{1}{m} \sum_{i=1}^m P(y^*|x^*, \theta^{(i)})$$

→ In practice (with DNN) it is common to use the whole training data for each of the m NNs. (random init., data shuffling typically ensures sufficient diversity.)

→ Note that in the above approximation, we sample from an approximate (empirical) posterior distribution \hat{P} using bootstrapping. ↗ not further explained

(multi-class) classification with NNs: NNs can also be used for classification.

For this we construct a NN with c outputs $F = [f_1, \dots, f_c]$, and normalize them into a probability distribution using softmax:

$$p = \text{softmax}(f), \text{ i.e. } p_i := \frac{\exp(f_i)}{\sum_{j=1}^c \exp(f_j)}, \quad P(y|x, \theta) = p_y. \quad (y: \text{class})$$

→ explicitly model the aleatoric uncertainty by injecting learnable (Gaussian) noise ϵ , and using $p = \text{softmax}(f + \epsilon)$.

Calibration: A model is well calibrated, if its confidence coincides with its accuracy across many predictions. Try adjusting the confidence estimation.

The Calibration error can be estimated by partitioning the test points into bins according to predicted confidence values (similar confidences together), then calculate the relative accuracy of pos. samples in each bin, and calculate the average confidences (weighted over bin sizes).

Reliability Diagrams: For each bin, plot the computed confidence against the frequency, and put all bins in one diagram. We want the histogram to move along the identity function.

→ Heuristics to improve model calibration:

- Histogram binning: Divide uncalibrated predictions of confidence \hat{P}_i into bins, and assign a new calibrated score $q_i = \text{freq}(B_m)$ to each bin.

- Isochoric regression: Find piecewise constant function $F = [f_1, \dots, f_m]$ that minimizes the bin-wise squared loss:

$$\min_{M, f, a} \sum_{m=1}^M \sum_{i=1}^n \mathbb{1}\{a_m \leq P(y_i=1|x_i) < a_{m+1}\} \cdot (f_m - y_i)^2 \quad (0 = a_1 \leq \dots \leq a_{M+1} = 1) \quad (f_1 \leq \dots \leq f_m)$$

Where f_1, \dots, f_m are the calibrated scores for each bin.

- Platt-scaling: adjust each value z_i of the output layer to $q_i = \sigma(a z_i + b)$, then learn a BLR to maximize the likelihood.

- Temperature scaling: Version of Platt-scaling with $b=0$ and $a=\frac{1}{T}$. T high: soft probabilities, T low: sharper probabilities.

- BNN usually improve calibration in practice, by themselves.

Active Learning: is about how to use uncertainty for deciding which data to collect. (which next sample is most useful?)

Mutual Information (information gain): Given r.v. X, Y we define

$$I(X, Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) \quad \text{the MI.} \quad (\rightarrow \text{symmetric})$$

The MI quantifies how much observing Y reduces uncertainty about X . Where $H(X|Y) = E_{Y \sim p(y)}[H(X|Y=y)] = E_{(x,y) \sim p(x,y)}[-\log(p(x|y))]$

is the conditional entropy, the expected remaining uncertainty in X , after observing Y . Note that $H(X|Y=y) \neq H(X|y)$, but it simply corresponds to a probabilistic update of our uncertainty in X , after observing the realization $y \sim Y$.

Further define the joint Entropy $H(X, Y) = E_{(x,y)}[-\log(p(x,y))]$

as the combined uncertainty about X and Y .

Note that $H(X, Y) = H(Y) + H(X|Y) = H(X) + H(Y|X)$.

Further $H(X|Y) \leq H(X)$, $I(X, Y) = H(X) + H(Y) - H(X, Y)$, $I(X, Y) \geq 0$.

"information never hurts principle"

Optimizing Mutual Information: Let f be a function which we are trying to approximate using Bayesian Learning (e.g. GP). Let D be set of points to potentially observe f at. We want to find $S \subseteq D$ which maximizes the information gain:

$$F(S) = H(f) + H(f|S) = I(f, S) \quad (= \frac{1}{2} \log(1 + \sigma^2_{f,S})) \quad \begin{matrix} \text{Example for } X \sim \mathcal{N}(0, 1) \\ y = x + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2) \end{matrix}$$

But: note that optimizing F is NP-hard.

A simple strategy is a Greedy algorithm: Let $S_t = \{x_1, \dots, x_t\}$, then choose

$$x_{t+1} = \arg \max_{x \in D} F(S_t \cup \{x\}) = \arg \max_{x \in D} \underbrace{F(S_t \cup \{x\}) - F(S_t)}_{=: F(x|S_t)} \quad \text{"marginal gain".}$$

We can define $I(f, Y_X | Y_{S_t}) = H(f|Y_{S_t}) - H(f|Y_X, Y_{S_t}) = F(S_t \cup \{x\}) - F(S_t)$.

For the example: $x_{t+1} = \arg \max_x \frac{1}{2} \log(1 + \frac{\sigma^2_f(x)}{\sigma^2_n}) = \arg \max_x \sigma^2_f(x)$.

→ "uncertainty sampling!"

→ The mutual information $F(S)$ is monotone submodular:

$$\forall x \in D, \forall A \subseteq B \subseteq D: F(A \cup \{x\}) - F(A) \geq F(B \cup \{x\}) - F(B) \quad \rightarrow \text{"diminishing returns property".}$$

→ Thus, the above Greedy algorithm provides a const.-factor

$$\text{approximation: } F(S_T) \geq (1 - \frac{1}{e}) \max_{S \subseteq D, |S| \leq T} F(S) \quad (1 - \frac{1}{e} \approx 0.63 !)$$

⇒ Uncertainty Sampling is near-optimal.

Uncertainty Sampling in heteroscedastic case: Fails to distinguish epistemic and aleatoric uncertainty. Because the most uncertain outcomes are not necessarily most informative. Maximizing mutual information yields: $x_{t+1} \in \arg \max_x \frac{\sigma^2_f(x)}{\sigma^2_n(x)}$ → epistemic uncertainty → aleatoric uncertainty.

Active Learning for Classification: In this setting a model produces a categorical distribution over the labels y_x for an input x . Uncertainty sampling corresponds to selecting samples that maximize the entropy of the predicted label:

$$x_{t+1} \in \arg \max_x H(Y|x_{1:t}, y_{1:t})$$

↓ → But note that this is often not a good approach, since the most uncertain label is not necessarily the most informative!

↓(AL for classification)

Bayesian active learning by disagreement (BALD): We can write

$$x_{t+1} = \operatorname{argmax}_x I(\theta, y_t | X_{1:t}, Y_{1:t}) = \operatorname{argmax}_x H(y_t | X_{1:t}, Y_{1:t}) - E_{\theta|X_{1:t}, Y_{1:t}} [H(y_t | \theta)] \quad (2)$$

(1) looks for points where the average prediction is not confident.

(2) look for points where models are confident in their predictions.

→ identifies points where models disagree about the label y_t .

Bayesian Optimization: Since data is costly, we want to efficiently improve our understanding of the data simultaneously to reaching certain goals (e.g. maximize f). For this, assume the following problem: Let $f^*: \mathcal{X} \rightarrow \mathbb{R}$, find $\operatorname{argmax}_{x \in \mathcal{X}} f^*(x)$

only able to obtain noisy observations $y_t = f^*(x_t) + \varepsilon_t$.

Further, these noisy observations are costly to obtain. To assume that similar alternatives yield similar results, we place a GP prior on f^* . → All following ideas generalize to other Bayesian Learning Problems.

Multi-armed-bandits (MAB) problem: We have k possible actions (arms) and want to maximize our reward, but no knowledge about the rewards distribution. Thus, we need to trade learning the reward distribution with following the most promising action.

Bayesian optimization is a version of MAB, with potentially infinite number of actions, but correlated reward.

Cumulative regret: The key performance metric in adaptively choosing new inputs to sample (online learning for BO) $y_t = f(x_t) + \varepsilon_t$ at, is called cumulative regret: $R_T = \sum_{t=1}^T (\max_x f(x) - f(x_t))$

The goal is to find algorithms, that achieve sublinear regret, i.e.

$$\lim_{T \rightarrow \infty} R_T / T = 0 \text{ thus } \max_t f(x_t) \rightarrow f(x^*).$$

→ We impose sublinearity, because with both algorithms (1) exploring forever, and (2) never exploring, the cumul. regret grows linearly.

→ sublinear regret requires balancing exploration and exploitation.

→ In Gaussian process bandit optimization we minimize R_T .

Upper Confidence Bound Sampling (GP-UCB): Uses a so-called acquisition function, to adaptively choose new inputs to sample. Here

$$x_{t+1} = \operatorname{argmax}_x M_t(x) + \beta_{\text{pri}} \cdot \sigma_t(x) \quad (\sigma_t(x) = \sqrt{k_t(x, x)}).$$

→ posterior mean and variance

β_t regulates how confident we are about our model f .

→ Optimism-based method: pick x with hope for best outcome.

→ this acquisition function is generally non-convex. (In low dim. use Lipschitz optimization; high-dim: gradient ascent with random init.)

Bayesian Regret of GP-UCB: Let $f \sim GP$, if we choose β_t "correctly",

$$\text{then } \frac{1}{T} \sum_{t=1}^T (f(x^*) - f(x_t)) = O^*(\sqrt{\frac{\log T}{T}}), \quad x_T = \max_{1 \leq t \leq T} I(f, y_s)$$

→ maximum information gain γ_T determines the regret.

(i.e. how quickly we can gain information)

→ Note that if the mutual information $I(f, y_s)$ is sublinear in T , then

R_T is sublinear, and we converge to the true optimum.

→ f^* assumed to be smooth (smaller class) ⇒ small information gain.

f^* assumed to be rough (larger class) ⇒ high information gain.

\downarrow (GP-UCB)

→ Information gain bounds for some common kernels:

- Linear: $\gamma_T = \mathcal{O}(d \cdot \log(T))$
 - Squared-exponential: $\gamma_T = \mathcal{O}((\log(T))^{d+1})$
 - Matérn with $\nu > \frac{1}{2}$: $\gamma_T = \mathcal{O}(T^{d/(2\nu+d)} \cdot (\log(T))^{2\nu/(2\nu+d)})$

→ This guarantees sublinear regret \Rightarrow convergence.

Frequentist regret for GP-UCB: Assume $f \in \mathcal{H}_K$ (reproducing kernel Hilbert space)

$$\text{then } \frac{1}{T} \sum_{t=1}^T (f(x^t) - f(x_t)) = \sigma^2 \left(\sqrt{\frac{\beta T}{T}} \right)$$

with $\beta_T = \sigma \left(\underbrace{\|f\|_K^2}_{\text{"complexity of } f\text{"}} + \gamma T \log(T)^3 \right)$ $\gamma T = \max_{1 \leq i \leq T} I(f_i, y_s)$

Thompson Sampling: Describes an acquisition function for Bayesian Optimization. At iteration t , we sample $\tilde{f} \sim P(F|X_{1:t}, Y_{1:t})$ and select $X_{t+1} = \operatorname{argmax}_{x \in D} \tilde{f}(x)$

→ randomness in realization of \tilde{f} is sufficient to trade exploration and exploitation.

→ We can establish regret bounds similar to GP-UCB.

(\rightarrow If not GP: sample from posterior distribution over models, then optimize the sample)

Kernel Hyperparameter Estimation: There are many different ways to learn the hyperparameters. In principle we could alternate learning hyperparameters and observation selection.

learning hyperparameters and observation selection.
But in practice, there is a specific danger of overfitting:
(1) Datasets are (in BO, AL) quite small, (2) Data points are selected
dependent on prior observations.

Some solutions to these problems are being Bayesian about the hyperparam. (i.e. place a hyperprior on them, and marginalize them out), or occasionally simply selecting some points at random.

(Finite) Markov Decision Process (MDP): is specified by a set of states $X = \{1, \dots, n\}$, actions $A = \{1, \dots, m\}$, transition probabilities (dynamics model) $P(X'|X, a) = P(X_{t+1} = x' | X_t = x, A_t = a)$,

and a reward function $r: X \times A \rightarrow \mathbb{R}$, mapping current state x and action a to some reward. r induces the sequence of rewards $R_t = r(X_t, A_t)$, $t \in \mathbb{N}$.

→ First assume we operate in a known environment (P, r known), which is fully observable (agent knows current state).

→ Goal: Choose actions which maximize the reward. Specifically we want to maximize the discounted payoff:

→ maximizes immediate and long-term reward.

Policy: is a function that maps each state $x \in X$ to a probability distribution over the actions, i.e.

$$\pi(a|x) = P(A_t=a|X_t=x) \quad (\text{deterministic case: } \pi: X \rightarrow A)$$

A policy π fully defines the behavior of an agent.

→ assume policies are stationary (i.e. do not change over time)

→ assume policies are stationary (i.e. do not change over time), → thus now view policies as deterministic.
 → See later: all optimal policies are deterministic!
 → E.g. policy iteration & Value Iteration (VPI) work with transition (simplification)

→ See later: all optimal policies are deterministic!
 → Every policy induces a Markov Chain $(X_t'')_{t \in \mathbb{N}}$ with transition probabilities $p''(x'|x) = P(X_{t+1}''=x'|X_t''=x) = P(x'|x, \pi(x)) = \sum_a \pi(a|x) \cdot P(x'|x, a)$

1

↑
7 deterministic

stochastic ↑

▼ (policies)

→ If the agent follows a fixed policy, then the evolution of the process can be fully explained by a Markov Chain.
→ It turns out we can find this policy quite efficiently.

State-action-value function: We try to quantify the effect of our initial state (and action) on our objective G_t (discounted payoff):

state-value function: $V_t^\pi(x) = E_\pi[G_t | X_t = x]$, and

→ lecture: $J(\pi) = E_\pi[G_0]$.

state-action-value function: $Q_t^\pi(x, a) = E_\pi[G_t | X_t = x, A_t = a]$

→ where we can also decide about the initial action

→ Note that $V^\pi(x) = Q^\pi(x, \pi(x))$ (if π deterministic).

→ Since we have stationary dynamics model, all quantities are independent of the start time t . → Thus write $V^\pi(x) := V_t^\pi(x)$, $Q^\pi(x, a) := Q_t^\pi(x, a)$. (wlog.)

Bellman Expectation Equations: V^π, Q^π fulfill the following equations:

$$V^\pi(x) = r(x, \pi(x)) + \gamma \cdot E_{x' | x, \pi(x)}[V^\pi(x')] = \sum_a \pi(a|x) (r(x, a) + \gamma \sum_{x'} P(x'|x, a) V^\pi(x'))$$

π stochastic
π deterministic

and

$$Q^\pi(x, a) = r(x, a) + \gamma \sum_{x'} P(x'|x, a) \cdot \sum_{a'} \pi(a'|x') \cdot Q^\pi(x', a')$$

$$\Rightarrow \text{deterministic case: } Q^\pi(x, \pi(x)) = r(x, \pi(x)) + \gamma \sum_{x'} P(x'|x, \pi(x)) Q^\pi(x', \pi(x))$$

→ The BEE gives us a recursive formula for the state-action-value function $V^\pi(x)$ (in deterministic case).

It follows that for a fixed policy π , we can find V^π by exactly!!

simply solving a system of linear equations:

$$V^\pi = r^\pi + \gamma P^\pi \cdot V^\pi \Leftrightarrow V^\pi = (I - \gamma P^\pi)^{-1} r^\pi$$

$$\text{where } V^\pi = \begin{pmatrix} V^\pi(1) \\ \vdots \\ V^\pi(n) \end{pmatrix}, \quad r^\pi = \begin{pmatrix} r(1, \pi(1)) \\ \vdots \\ r(n, \pi(n)) \end{pmatrix}, \quad P^\pi = \begin{pmatrix} P(1|1, \pi(1)) & \cdots & P(n|1, \pi(1)) \\ \vdots & \ddots & \vdots \\ P(1|n, \pi(n)) & \cdots & P(n|n, \pi(n)) \end{pmatrix}$$

approximately

Fixed-point Iteration: Can be used to solve the above system of linear equations for V^π . We can show that V^π is the unique fixed-point of the affine mapping: $B^\pi V = r^\pi + \gamma P^\pi V$.

Use the algorithm: (1) initialize V^π (randomly), (2) update the value

$$V^\pi \leftarrow B^\pi V^\pi = r^\pi + \gamma P^\pi V^\pi \text{ (until convergence).}$$

→ Fixed-point iteration converges to V^π exponentially quickly.

Policy Optimization: We denote our optimal policy to maximize the discounted payoff: $\pi^* = \underset{\pi}{\operatorname{argmax}} E_\pi[G_0]$ → irresp. of initial state!

Alternatively say we want the maximal element π^* of the partial ordering $\pi \geq \pi' \Leftrightarrow V^\pi(x) \geq V^{\pi'}(x) \forall x \in X$.

→ all optimal policies have identical value function, thus write

$$V^*(x) = V^{\pi^*}(x) = \max_{\pi} V^\pi(x), \quad Q^*(x) = Q^{\pi^*}(x) = \max_{\pi} Q^\pi(x).$$

Greedy Policies: The Greedy policy w.r.t. a state-action value function Q (based on any π') is given by

$$\pi_Q(x) = \underset{a \in A}{\operatorname{argmax}} Q(x, a) \rightarrow \text{choose action which maximizes the discounted payoff.}$$

Analogously for a state value function V :

$$\pi_V(x) = \underset{a \in A}{\operatorname{argmax}} r(x, a) + \gamma \sum_{x'} P(x'|x, a) \cdot V(x')$$

→ inserting Bellmann expectation equality for $V(x)$.

↓ (greedy policies)

→ The greedy policy π^* induces a new value function V^{π^*} .
 With this new value function, we can then again obtain a new
 greedy policy, inducing then again a new new value function ...
 → This shows a cyclic dependency.

Bellman Theorem: It turns out, the optimal policy π^* is a fixed-point
 of the above described cyclic dependency. In particular
 π^* is optimal \Leftrightarrow if π^* is optimal w.r.t its own value function

Note: $V^*(x) = \max_a Q^*(x, a) = \max_a r(x, a) + \gamma \sum_{x'} P(x'|x, a) V^*(x')$ → Bellmann optimality equation

→ This Bellmann optimality equation is non-linear (→ generally no closed-form solution). Thus typically use iterative methods to solve.
 → This proves: there always exists an optimal policy π^* . And π^* is deterministic and stationary

Policy Iteration: is a way of computing the optimal policy.

- (1) Start with an arbitrary initial policy.
 - (2) Use the Bellman expectation equation (solve linear system of equations), to find the value function of that policy,
 - (3) then choose the greedy policy (w.r.t. the new value function) as the next iterate of the algorithm (→ then go back to (1)).
- is guaranteed to monotonically improve (i.e. $V^{\pi_{t+1}}(x) \geq V^{\pi_t}(x) \forall t, x$)
 → converges to optimal policy π^* in polynomial time $\mathcal{O}(n^2 m / (1-\gamma))$.

Value Iteration: is a way of computing the optimal policy by fixed-point iteration.

For this we regard V^* as the fixed point of the Bellmann-optimality equation $V^*(x) = \max_a r(x, a) + \gamma \sum_{x'} P(x'|x, a) V^*(x)$.

Algorithm:

- (1) Initialize $v_0(x) = \max_a r(x, a)$ for all $x \in X$,

(2) Iteratively (int), let $Q_t(x, a) = r(x, a) + \gamma \sum_{x'} P(x'|x, a) v_{t-1}(x)$

then for each x , let: $v_t(x) = \max_a Q_t(x, a)$ → for some previously chosen ε

(3) Break if $\|v_t - v_{t-1}\|_\infty = \max_x |v_t(x) - v_{t-1}(x)| \leq \varepsilon$.

(4) Choose greedy policy w.r.t. v_t .

→ converges to an ε -optimal solution in polynomial time.

→ but in general does not converge to optimal solution.

→ In sparse MDP, one iteration of Value Iteration has constant time. → i.e. sparse connections of state space.

Partially Observable Markov Decision Processes (POMDP): Let us assume the partially observed setting (i.e. agent can only access noisy observations y_t of its state x_t).

Define a POMDP by states X , actions A , transition probabilities $P(x'|x, a)$, and reward function $r: X \times A \rightarrow \mathbb{R}$. Further define observations y , and observation probabilities $\alpha(y|x) = P(y_t = y | x_t = x)$. POMDPs are controlled hidden Markov models.

→ generally a very powerful model, but typically extremely intractable.

POMDP can be reduced to an MDP with enlarged state space.

For this consider an MDP with states $b_t(x) = P(X_t = x | y_{1:t}, a_{1:t})$ "beliefs".

i.e. each state is a probability distribution over states of POMDP.

We keep track of how our beliefs change over time (i.e. Bayesian Filtering). Given a prior belief b_t , and action a_t taken, and observation y_t , the belief state can be updated as follows:

$$b_{t+1}(x) = P(X_{t+1} = x | y_{1:t+1}, a_{1:t}) = \frac{1}{Z} \cdot \alpha(y_{t+1}|x) \sum_{x' \in X} P(x'|x, a_t) \cdot b_t(x')$$

↓(POMDPs)

→ can use MDP methods to solve (theoretically.)

This leads us to a belief state MDP, defined by

(1) belief space $B = \Delta^X$ (hidden states X), (2) actions A , (3) transition probabilities $T(b'|b,a) = P(B_{t+1} = b' | B_t = b, A_t = a)$, and (4) rewards $r(a,b) = E_{x \sim b}[r(x,a)] = \sum_{x \in X} b(x) r(x,a)$

→ For finite horizon T (in discounted payoff), set of reachable belief states is finite (but exponential in T) → could use dynamic programming.

→ most belief states of POMDPs are never reached. Thus it is common to discretize the belief space (e.g. sampling, dimensionality reduction).

↳ alternative approach: Policy Gradients (see later).

Reinforcement Learning: is concerned with probabilistic planning in unknown environments. These environments can be modelled using MDP's, but with unknown dynamics p and rewards r .

We focus on the fully observed setting (agent knows current state), and, for now, on small state and action spaces (i.e. tabular setting).

→ exploration-exploitation dilemma applies (like in Bayesian opt.)

↳ Bayesian Optimization = RL with fixed state and continuous action space.

→ Data is not i.i.d! But depends on played actions. → compared to supervised learning.

The collected data is modelled by so-called trajectories

$\tau = (x_0, a_0, r_0, x_1, a_1, r_1, x_2, \dots)$. There are two common settings

• episodic: Agent learns over multiple "training" episodes i . Each episode results in a new trajectory τ_i , after which the environment resets.

• non-episodic/online: Agent learns "online", single trajectory.

On-policy RL methods: is where the agent has full control over which actions to pick. Can experiment with exploration and exploitation.

Off-policy RL methods: Agent does not (always) have control over actions (only observational data). → more sample efficient. (than on-policy)

→ also that all observational actions are wrt. the same policy (doesn't change)

Model-based RL approach: learn underlying MDP, by estimating the dynamics model p and reward r . Then use these estimates to perform planning (i.e. policy optimization).

Model-free RL approach: Estimate the value function directly. → actor critic methods

Learn underlying MDP: One way to do this, is to use MLE/MAP.

Such estimates are: $\hat{p}(x_{t+1}|x_t, a) = \frac{N(x_{t+1}|x_t, a)}{N(a|x)}$ $\tau = (x_0, a_0, r_0, x_1, \dots)$

where $N(x_{t+1}|x_t, a)$: # transitions from x_t to x_{t+1} with action a .

$N(a|x)$: # transitions that start in x_t , and play a .

and the rewards: $\hat{r}(x, a) = \frac{1}{N(a|x)} \sum_{t=0, x_t=x, a_t=a}^{\infty} R_t$

→ both estimates are unbiased, and correspond to the sample mean.

→ only accurate if each pair (x, a) is visited many times!

on-policy, model-based

ϵ -greedy algorithm: is a way of choosing our actions to best estimate p, r .

At each time step t with probability ϵ_t we pick an action uniformly at random, and with probability $1 - \epsilon_t$ we pick the best action under our current model.

If $\sum_t \epsilon_t = \infty$ and $\sum_t \epsilon_t^2 < \infty$, then we can converge to the optimal policy with probability 1. (using Monte Carlo control)

→ often performing fairly well

→ doesn't quickly eliminate suboptimal actions.

Rmax-Algorithm: is a way of estimating our dynamics model P , and rewards r . Idea: Optimism in the face of uncertainty.
 This means, if $r(x,a)$ is unknown, we set $\hat{r}(x,a) = R_{\max}$ (i.e. the maximum possible reward). Similarly, if $p(x'|x,a)$ is unknown, we imagine an additional "fairy-tale" state x^* , and set $\hat{P}(x^*|x,a) = 1$.
 x^* fulfills: $\hat{P}(x^*|x^*,a) = 1 \quad \forall a, \quad \hat{r}(x^*,a) = R_{\max} : \quad \boxed{x} \xrightarrow{R_{\max}} \boxed{x^*} \circlearrowleft R_{\max}$.

→ If in doubt, the agent believes actions to lead to a fairy tale state x^* with maximal rewards. → Encourages exploration of unknown states.
 In particular, the Rmax-algorithm first computes the optimal policy π^* from \hat{P} and \hat{r} , then executes π^* for some number of steps, then updates \hat{P} and \hat{r} , then recomputes the optimal policy from \hat{P} and \hat{r} , and then repeats the process.
 (→ For each state we explore all actions, and then always choose the best) → for all x in parallel

- For some T , every T timesteps, with high probability R_{\max} either obtains near optimal reward, or visits at least one unknown state-action pair. (T is related to mixing time of MC).
- With probability $1-\delta$, R_{\max} reaches an ϵ -optimal policy in time polynomially in $|X|, |A|, T, 1/\epsilon, \log(1/\delta)$, and R_{\max} -value.
 This can be shown using the Hoeffding bound: Z_1, \dots, Z_n i.i.d. bounded in $[0,C]$ with mean M fullfill:

$$P(|M - \frac{1}{n} \sum_{i=1}^n Z_i| > \epsilon) \leq 2 \exp(-2n\epsilon^2/C^2)$$

→ Rmax performs remarkably well in tabular setting.

Challenges of Model-based approaches:

Computational limitations: $\hat{P}(x'|x,a)$ and $\hat{r}(x,a)$ need to be stored in a table with $\mathcal{O}(n^2m)$ entries. → quickly unmanageable.

Also we need to "solve" the learned MDP, to obtain the optimal policy (many times). → Rmax: $\mathcal{O}(nm)$ times.

On-Policy Value estimation: In model-free RL approaches, we estimate the value function V^π directly. Remember

$$\begin{aligned} V_{\text{true}}^\pi(x) &= r(x, \pi(x)) + \gamma \sum_{x'} p(x'|x, \pi(x)) \cdot V_{\text{old}}^\pi(x') = \mathbb{E}_{x' \sim P} [R + \gamma V^\pi(x') | x, \pi(x)] \\ &\approx r + \gamma V^\pi(x') \quad \text{for observed } (x, a, r, x') \end{aligned}$$

This approx. corresponds to a one-sample Monte Carlo approx. → do this repeatedly i.e. each iteration.

But since we don't know the true $V^\pi(x)$, we use a bootstrapping estimate $V_{\text{old}}^\pi(x)$. (bootstrapping: estimate true quantity with an empirical quantity.)

→ bootstrapping is biased. We have high variance, because single sample.

Temporal Difference (TD-) Learning: is on-policy, model free approach where we reduce the variance of the bootstrapping estimate above. Particularly, first initialize V^π arbitrarily. Then for each iteration t first, follow policy π to obtain transition (x, r, a, x') , then update

$$V^\pi(x) \leftarrow (1-\alpha_t) V^\pi(x) + \alpha_t (r + \gamma V^\pi(x')) \quad \text{or equivalently}$$

$$[V^\pi(x) \leftarrow V^\pi(x) + \alpha_t (r + \gamma V^\pi(x') - V^\pi(x))].$$

→ on-policy TD-L for Q^π is called SARSA (same guarantees)

(Thm)
 If $\sum_t \alpha_t = \infty$, $\sum_t \alpha_t^2 < \infty$ and all state-action pairs are chosen infinitely often, then \hat{V}^π converges to the true V^π with probability 1.

Off-policy version: The actions a don't need to be chosen via π . Thus off-policy is also possible. The Thm above also holds for Q^π :

$$\text{update: } Q^\pi(x,a) \leftarrow (1-\alpha_t) Q^\pi(x,a) + \alpha_t (r + \gamma Q^\pi(x', \pi(x'))).$$

Q-Learning: is an off-policy, model free approach, where we approximate the value function Q^* of the optimal policy π^* directly. For this, remember from the Bellman Theorem:

$$\pi^*(x) = \arg \max_{a \in A} Q^*(x, a), \text{ and}$$

and that $Q^*(x, a) = r(x, a) + \gamma \sum_{x' \in X} p(x'|x, a) \cdot \max_{a' \in A} Q^*(x', a')$

$$\approx r + \gamma \max_{a' \in A} Q^*(x', a')$$

→ one-sample Monte Carlo approximation

Using a bootstrapping estimate for this approximation results in the Q-Learning algorithm: (1) Initialize Q^* arbitrarily.

(2) For each iteration t , observe the transition (x, a, r, x') , then update:

$$Q^*(x, a) \leftarrow (1-\alpha_t) \cdot Q^*(x, a) + \alpha_t (r + \gamma \cdot \max_{a' \in A} Q^*(x', a'))$$

If $\sum_t \alpha_t = \infty$, $\sum_t \alpha_t^2 < \infty$, and all state-action pairs are chosen infinitely often, then Q^* converges to optimal Q^* with probability 1.
→ with probability at least $1-\delta$, Q-learning converges to an ϵ -optimal policy in time polynomial in $\log(|X|)$, $\log(|A|)$, $1/\epsilon$, $\log(1/\delta)$.

On-policy version: Optimistic Q-Learning. (is similar to Rmax).

(1) Initialize $Q^*(x, a) = V_{\max} \cdot \prod_{t=1}^{T_{\text{init}}} (1-\alpha_t)^{-1}$,

(2) For each iteration t choose action $a_t = \arg \max_{a \in A} Q^*(x, a)$, and observe (x, a, r, x') , then update

$$Q^*(x, a) \leftarrow (1-\alpha_t) Q^*(x, a) + \alpha_t (r + \gamma \cdot \max_{a' \in A} Q^*(x', a'))$$

where $V_{\max} = R_{\max} / (1-\gamma) \geq \max_{x, a} Q^*(x, a)$, T_{init} some initialization time.

→ with prob. at least $1-\delta$, optimistic Q-L conv. to ϵ -optim. policy in time poly. in $|X|$, $|A|$, $1/\epsilon$, $\log(1/\delta)$.

→ Q-Learning requires $\Theta(nm)$ memory, and during each iteration $\Theta(m)$.

Challenges of Model-free Approaches: Model-free Q-learning takes

polynomial time in $|X|$ and $|A|$. But in non-tabular settings (i.e. large state/action spaces) this is unacceptable. Because

often domains are continuous, e.g. when modelling beliefs about states in a partially observable environment (POMDPs)

→ Thus we want to learn/approximate value functions (regression). (parameterization)

(tabular)

TD-Learning as SGD: We can reinterpret model-free methods (such as TD-1/Q-Learning) as solving an optimization problem, where each iteration corresponds to a single gradient update.

Let us parameterize our estimate function of V^{π} , with param. θ .

The TD-Learning update rule: $V^{\pi}(x) \leftarrow (1-\alpha_t) V^{\pi}(x) + \alpha_t (r + \gamma V^{\pi}(x'))$

can equivalently be viewed as stochastic (semi-) gradient descent on the squared loss $L_2(\theta; x, x', r) = \frac{1}{2} (V(x; \theta) - r - \gamma V(x'; \theta_{\text{old}}))^2$

→ the term inside the $(\cdot)^2$, we call "TD-error" = δ

Update via SGD: $V \leftarrow V - \alpha_t \delta = \hat{V}^{\pi}(x; \theta_{\text{old}}) - \alpha_t \cdot (\hat{V}^{\pi}(x; \theta_{\text{old}}) - r - \gamma \hat{V}^{\pi}(x'; \theta_{\text{old}}))$

$$= (1-\alpha_t) \hat{V}^{\pi}(x; \theta_{\text{old}}) + \alpha_t \cdot (r + \gamma \hat{V}^{\pi}(x'; \theta_{\text{old}}))$$

→ bootstrapping: we use θ_{old} value estimates as labels (and opt. wrt. them), $\hat{V}^{\pi}(x'; \theta_{\text{old}})$ not wrt. θ . Note that the optimization target $r + \gamma \hat{V}^{\pi}(x'; \theta_{\text{old}})$ moves between iterations. ↳ stability issues.

→ The same insight applies to learning the (optimal) action-value function Q .

→ This shows a path towards parametric function approximation.

Parametric Value function approximation: To scale to large state spaces, we approximate the value function using a parameterized model $V(x; \theta)$ or $Q(x, a; \theta)$. This is a strict generalization of the tabular setting. A straightforward approach is to use a linear function approximation with hand-designed feature map ϕ : $Q(x, a; \theta) = \theta^T \phi(x, a)$.

Alternatively it is common to use deep neural networks. This approach is known as deep reinforcement learning.

Q-Learning with function approximation: we can generalize Q-Learning from the tabular setting. After observing a transition (x, a, r, x') , we update the parameters θ via the gradient of:

$$L_2(\theta; x, a, r, x') = \frac{1}{2} (\hat{Q}^*(x, a; \theta) - r - \gamma \max_{a' \in A} Q^*(x', a'; \theta_{old}))^2 \quad \text{"Bellman error"}$$

Thus $\theta \leftarrow \theta - \alpha_t \cdot \delta \cdot \nabla_\theta Q(x, a; \theta)$ with $\delta := Q(x, a; \theta) - r - \gamma \max_{a'} Q(x', a'; \theta)$.

→ In tabular setting, this is identical to Q-Learning.

→ Typically rather slow algorithm. (In practice convergence if function class "rich enough")
→ moving gradient target (with each iteration) → stability issues.

(Double) Deep Q-Networks (DQN/DDQN): DQN aims to "stabilize" the optimization targets in the above gradient descend. For this it updates the NN used for the approximate bootstrapping estimate infrequently, to maintain a constant optimization target across multiple episodes.

Exact implementations of DQN vary. One approach is to clone the NN , and maintain one changing "online network" for most recent estimate of Q (parameterized by θ), and one fixed "target network" used as the optimization target (parameterized by θ_{old}), which is updated infrequently. (i.e. θ_{old} is only updated to θ after some amount of observations).

Thus in DQN does gradient descent according to the Loss:

$$L_{DQN}(\theta; D) = \frac{1}{2} \sum_{(x, a, r, x') \in D} (r + \max_{a' \in A} Q^*(x', a'; \theta_{old}) - Q^*(x, a; \theta))^2 \quad \text{→ } a' \text{ is chosen wrt. } \theta_{old} !!$$

→ Maximization Bias: The estimates of the true Q^* are noisy. But since we maximize noisy values, we receive a biased estimate.

DDQN addresses the problem of maximization bias by choosing the optimal action w.r.t. the new network (instead of θ_{old} in DQN).

$$\text{Loss: } L_{DDQN}(\theta; D) = \frac{1}{2} \sum_{(x, a, r, x') \in D} (r + \gamma \cdot \underbrace{Q^*(x', a^*(x'; \theta); \theta_{old})}_{\text{opt chosen wrt. new network } \theta} - Q^*(x, a; \theta))^2$$

$$\text{where } a^*(x'; \theta) = \arg \max_{a' \in A} Q^*(x', a'; \theta) \quad \text{→ opt chosen wrt. new network } \theta. \quad (\text{not } \theta_{old})$$

Thus, online update: $\theta \leftarrow \theta + \alpha_t (r + \gamma \cdot Q^*(x', a^*(x'; \theta); \theta_{old}) - Q^*(x, a; \theta)) \nabla_\theta Q^*(x, a; \theta)$

after observing a single transition (x, a, r, x') , $a^*(x'; \theta)$ is treated as a const. when differentiating. θ_{old} is updated to θ after $|D|$ -many transitions observed.

Policy Approximation: Another way of scaling to large (continuous) action spaces, is to directly learn an approximate parameterized policy $\pi^*(x) \approx \pi(x; \theta) = \pi_\theta(x)$ θ : parameters

using e.g. NN. Policy search and policy gradient are such methods.

Note: $\pi_\theta(x)$ is a probability distribution over A . Policy gradient methods fundamentally rely on randomized policies for exploration.

Policy Value Function: For approximating the optimal policy, we will focus on the discounted payoff through the policy value function

$$J(\theta) = E_{\pi_\theta} \left[\sum_{t=0}^{\infty} \gamma^t \cdot r(x_t, \pi_\theta(x_t)) \right], \quad \text{= } G_0 \text{ (discounted payoff)} \rightarrow \begin{array}{l} \text{since } \pi_\theta(x) \text{ is a} \\ \text{probability distribution,} \\ \text{we can only} \\ \text{compute an expected value.} \end{array}$$

the expected discounted payoff of π_θ .

Bounded Variant: $J_T(\theta) = E_{\pi_\theta} \left[\sum_{t=0}^{T-1} \gamma^t \cdot r(x_t, \pi_\theta(x_t)) \right] = E_{\pi_\theta}[G_{0:T}]$

The general idea is to find the optimal policy by maximizing $J(\theta)$
i.e. to find $\theta^* = \arg \max J(\theta)$.

Such methods we call Policy search methods.

Note that we can approximate $J(\theta)$ well in episodic tasks (i.e. agent can be "reset"), using "rollouts". Let $\mathcal{T}^{(1)}, \dots, \mathcal{T}^{(m)} \sim \pi_\theta$ trajectories
 $(\mathcal{T}^{(i)}) = (x_0^{(i)}, a_0^{(i)}, r_0^{(i)}, x_1^{(i)}, \dots)$ each of length T . For each trajectory we can compute the i -th (finite horizon) discounted payoff: $G_{0:T}^{(i)} = \sum_{t=0}^{T-1} \gamma^t r_t^{(i)}$. is here seen as induced distribution over trajectories.

Then we can use Monte Carlo Sampling to approximate $J_T(\theta)$.

Thus $J(\theta) \approx J_T(\theta) \approx \frac{1}{m} \sum_{i=1}^m G_{0:T}^{(i)}$

→ on-policy (fundamentally)

Policy Gradients: To maximize $J(\theta)$, we want to compute its gradients. We can show that the following holds:

$$\nabla_\theta J(\theta) = \nabla_\theta J_T(\theta) = \nabla_\theta E_{\mathcal{T} \sim \pi_\theta} [r(\mathcal{T})] = E_{\mathcal{T} \sim \pi_\theta} [r(\mathcal{T}) \nabla_\theta \log(\pi_\theta(\mathcal{T}))]$$

where $r(\mathcal{T}) := \sum_{t=0}^T \gamma^t \cdot r(x_t, a_t)$ (also: $\nabla_\theta \log(\pi_\theta(\mathcal{T})) = \frac{\nabla_\theta \pi_\theta(\mathcal{T})}{\pi_\theta(\mathcal{T})}$)

$\hookrightarrow \nabla_\theta \pi_\theta(x) = \pi_\theta(x) \cdot \nabla_\theta \log(\pi_\theta(x))$ "score gradient trick"

For estimating this gradient, we need to evaluate $\nabla_\theta \log(\pi_\theta(\mathcal{T}))$:

Note that $\nabla_\theta \log(\pi_\theta(\mathcal{T})) = \sum_{t=0}^{T-1} \nabla_\theta \log(\pi_\theta(a_t, x_t))$

using $\pi_\theta(\mathcal{T}) = p(x_0) \cdot \prod_{t=0}^{T-1} \pi_\theta(a_t | x_t) \cdot P(x_{t+1} | x_t, a_t)$ → exploiting MDP structure.

Thus $\nabla_\theta J(\theta) \approx E_{\mathcal{T} \sim \pi_\theta} [r(\mathcal{T}) \cdot \sum_{t=0}^{T-1} \nabla_\theta \log(\pi_\theta(a_t, x_t))]$

$$\approx \frac{1}{m} \cdot \sum_{i=1}^m G_{0:T}^{(i)} \sum_{t=0}^{T-1} \nabla_\theta \log(\pi_\theta(a_t^{(i)}, x_t^{(i)})) \rightarrow \begin{array}{l} \text{Monte} \\ \text{Carlo Sampling.} \\ \rightarrow \text{episodic setting} \end{array}$$

→ Estimates are unbiased, but typically have large variance.

→ optimized policy can be randomized.

Baselines: We can reduce variance of policy gradient estimates, by using so-called baselines b . Note that

$$\forall b \in \mathbb{R}: E_{\mathcal{T} \sim \pi_\theta} [r(\mathcal{T}) \nabla_\theta \log(\pi_\theta(\mathcal{T}))] = E_{\mathcal{T} \sim \pi_\theta} [(r(\mathcal{T}) - b) \cdot \nabla_\theta \log(\pi_\theta(\mathcal{T}))]$$

a common baseline is $b(\mathcal{T}_{0:T-1}) = \sum_{k=0}^{T-1} \gamma^k r_k$, also $G_0 - b(\mathcal{T}_{0:T-1}) = \gamma^T G_{t:T}$. ↑ inserted term with b simplifies to 0.

Thus yielding the gradient estimator:

$$\nabla_\theta J(\theta) \approx \nabla_\theta J_T(\theta) = E_{\mathcal{T} \sim \pi_\theta} \left[\sum_{t=0}^{T-1} \gamma^t G_{t:T} \cdot \nabla_\theta \log(\pi_\theta(a_t | x_t)) \right] \quad \text{non-policy}$$

↑ "downstream return"

"reward to go" after action a_t

REINFORCE algorithm: performs SGD on the above reduced-variance policy gradient. Specifically: (1) initialize policy weights θ , then

(2) generate an episode (i.e. rollout) and obtain on trajectory \mathcal{T} .

(2.1) For each iteration $t=0$ to $t=T-1$ set $G_{t:T}$ to the return from iteration t , and perform the update

$$\theta \leftarrow \theta + \eta \cdot \gamma^t G_{t:T} \nabla_\theta \log(\pi_\theta(a_t | x_t))$$

(3) After all iterations, go back to (2), and generate a new trajectory.



↓ (REINFORCE)

The variance of REINFORCE can be further reduced by subtracting from $G_{t:T}$: $\nabla_\theta J(\theta) \approx E_{G \sim \pi_\theta} \left[\sum_{t=0}^T \gamma t (G_{t:T} - b_t) \nabla_\theta \log(\pi_\theta(a_t|x_t)) \right]$

where b_t is the mean reward to go: $b_t = \frac{1}{T-t} \sum_{t=0}^{T-1} G_{t:T}$

→ main advantage of REINFORCE: can be used in continuous action spaces.

→ BUT: not guaranteed to find optimal policy. Even in small domains, REINFORCE can get stuck in local optima

→ policy gradient methods are slow. (large variance in gradient estimates → small steps, and need many MC rollouts)

→ We cannot reuse data, since policy gradient methods are fundamentally on-policy (since estimation relies on policy not changing) → Samples rollouts using the current estimate of the policy.

Policy Gradient Theorem: We can reinterpret the Policy Gradient in terms of the state-action value function Q , such that

$$\begin{aligned} \nabla_\theta J(\theta) &= E_{G \sim \pi_\theta} \left[\sum_{t=0}^\infty \gamma^t \cdot Q(x_t, a_t) \cdot \nabla_\theta \log(\pi_\theta(a_t|x_t)) \right] \\ &= \int p_\theta(x) \cdot E_{a \sim \pi_\theta(x)} [Q(x, a) \nabla_\theta \log(\pi_\theta(a|x))] dx \\ &=: E_{(x,a) \sim \pi_\theta} [Q(x, a) \cdot \nabla_\theta \log(\pi_\theta(a|x))] \end{aligned}$$

where $p_\theta := \sum_{t=0}^\infty \gamma^t p(x_t=x)$ is the (non-normalized, discounted) state occupancy measure.

→ Since we do not know the true Q , we will approximate $Q_\phi(x, a)$ using parameters ϕ (e.g. bootstrapping).

On-policy Actor Critics: (online setting) Actor Critic methods consist of a parameterization of the policy π_θ ("actor"), and of the value function Q_ϕ ("critic"). In DRL we use NN to parameterize both actor and critic. → allows scaling to large state/action spaces.

One example in the online setting is to use TD-Learning to learn the critic and SGD for the actor. For some initial ϕ, θ , we do for each iteration: (1) use π_θ to obtain (x, a, r, x') .

(2) set $\delta = r + \gamma Q_\phi(x', \pi_\theta(x')) - Q_\phi(x, a)$, (3) then do the updates

$$\begin{aligned} \theta &\leftarrow \theta + \eta_1 \gamma t \cdot Q_\phi(x, a) \cdot \nabla_\theta \log(\pi_\theta(a|x)) && \text{→ actor update} \\ \phi &\leftarrow \phi + \eta_2 \cdot \delta \cdot \nabla_\phi Q_\phi(x, a) && \text{→ critic update} \end{aligned}$$

we disregard dependence of Q on θ .

→ Note: the actor is not guaranteed to improve, due to bias in estimates.

To further reduce variance of the above estimates, we can introduce a baseline into estimate of our actor:

$$\theta \leftarrow \theta + \eta_1 \gamma t (Q_\phi(x, a) - V_\phi(x)) \cdot \nabla_\theta \log(\pi_\theta(a|x))$$

$$A(x, a) = Q(x, a) - V(x)$$

where $A_\phi(x, a) = Q_\phi(x, a) - V_\phi(x)$ is an estimate of the advantage function.

⇒ We call this algorithm A2C (advantage actor critic). → easier to estimate

→ Actor Critic implementations often use advantage function. than Q .

Algorithms like GAE (generalized advantage estimation) address the bias-variance tradeoff by blending the estimates of REINFORCE (high variance, but unbiased) with using a bootstrapped Q -function (smaller variance, but biased) for the gradient estimates.

→ Actor Critic methods usually rely on randomization of the policy, for exploration. This is often insufficient. → use ϵ -greedy policy.

→ low sample efficiency. On-policy actor critic methods need extremely large number of interactions, because they cannot reuse past data. → TRPO / PPO address this.

→ because policy quickly collapses into deterministic policy.

TRPO (trust region Policy optimization): is a modern variant of Actor-Critic methods, which slightly improves sample efficiency.
During each iteration we use a fixed critic to optimize the policy, particularly: $\theta_{k+1} \leftarrow \underset{\theta}{\operatorname{argmax}} J(\theta) \text{ s.t. } KL(\theta || \theta_k) \leq \delta$

$$\text{where } J(\theta_k, \theta) = E_{(x,a) \sim \pi_{\theta_k}} \left[\frac{\pi_{\theta}(a|x)}{\pi_{\theta_k}(a|x)} \cdot A_{\theta_k}(x,a) \right] \quad \xrightarrow{\text{advantage function}}$$

→ can reuse data from rollouts within the same iteration, as long as they can still be "trusted" → "somewhat" off-policy → fundamentally still on-policy.

PPO (proximal policy optimization) is a variant of TRPO, which works well in practice. It controls importance of weights directly, instead of imposing a constraint on KL-divergence. → GPT uses PPO.

Off-policy Actor Critics: In many applications sample efficiency is crucial, which is why we look at off-policy methods, where past data can be reused. Such algorithms use reparameterization gradient estimates (like in Variational Inference), instead of score gradient estimators. These methods are closely related to value iteration, to learn the optimal value function, to then find the optimal policy.

For now assume the policy π is deterministic.

$$\text{Remember the DQN loss: } L_{\text{DQN}}(\theta; D) = \frac{1}{2} \sum_{(x,a,r,x') \in D} (r + \gamma \max_{a' \in A} Q_{\theta_{\text{old}}}^*(x', a') - Q_{\theta}(x, a))^2$$

to solve the maximization problem of Q^* if the action space A is very large, we can replace the maximum with a parameterized policy:

$$L(\varphi) = \frac{1}{2} \sum_{(x,a,r,x') \in D} (r + \gamma \cdot Q_{\varphi_{\text{old}}}(x', \underline{\pi}_{\theta}(x')) - Q_{\varphi}(x, a))^2$$

By the Bellman Theorem, we want π_{θ} to approximate the greedy policy:

$$\pi_{\theta}(x) \approx \pi^*(x) = \underset{a \in A}{\operatorname{argmax}} Q_{\varphi}^*(x, a).$$

And if we allow "rich enough" policies, then this is equivalent to:

$$\theta^* \in \underset{\theta}{\operatorname{argmax}} E_{x \sim \mu} [Q(x, \pi_{\theta}(x))] \quad \text{where } \mu(x) > 0 : \text{"explores all states!"}$$

Then, we use a bootstrapping estimate of Q^* , and thus neglect the dependence of the critic Q^* on the policy (actor) parameters θ . We can use chain rule: $D_{\theta} Q_{\varphi}^*(x, \pi_{\theta}(x)) = D_{\theta} Q_{\varphi}^*(x, a) | a = \pi_{\theta}(x)$.

Exploration: Usually policy gradient techniques rely on randomness of the policy. But here we assume deterministic policies.

→ DDPG (deep deterministic policy gradients) addresses this problem, by injecting Gaussian noise into actions selected by π_{θ} . → "Gaussian noise dithering"
This is essentially equivalent to using Q-Learning, but replace max with $\max_{a \sim \pi_{\theta}(x)}$. Particularly: (1) Initialize $\theta, \varphi, D=\emptyset$; replay buffer, $\theta_{\text{old}} = \theta, \varphi_{\text{old}} = \varphi$.
repeat: (2) Observe state x , pick action $a = \pi_{\theta}(x) + \epsilon, \epsilon \sim \mathcal{N}(0, \lambda I)$.

(3) execute and observe (x, a, r, x') , store this in D .

(4) If we have made enough observations, do

(4.1) Sample mini batch $B \subseteq D$, for each observation (in batch)

compute $y = r + \gamma \cdot Q_{\varphi_{\text{old}}}(x', \pi_{\theta_{\text{old}}}(x'))$, and update

$$\varphi \leftarrow \varphi - \eta \nabla_{\varphi} \frac{1}{|B|} \sum_B (Q_{\varphi}(x, a) - y)^2 \quad \xrightarrow{\text{critic update}}$$

$$\theta \leftarrow \theta + \eta \nabla_{\theta} \frac{1}{|B|} \sum_B Q_{\varphi}(x, \pi_{\theta}(x)) \quad \xrightarrow{\text{actor update}}$$

$$\theta_{\text{old}} \leftarrow (1-\rho) \theta_{\text{old}} + \rho \cdot \theta \quad \text{AND} \quad \varphi_{\text{old}} \leftarrow (1-\rho) \varphi_{\text{old}} + \rho \cdot \varphi$$

→ TD3 (Twin Delayed DDPG) extends DDPG by using two critic networks, for predicting maximum action, and for evaluating the policy. This addresses maximization bias. It also increases stability, by delaying updates to the actor network.

Randomized policies in off-policy Actor Critics: To ensure exploration, instead of injecting noise (e.g. DDPG), we can directly allow randomized policies.

Looking at DDPG, instead of using the critic update (for Q), we can obtain unbiased gradient estimates, by sampling $a \sim \pi_{\text{old}}(x')$. Further suppose the policy is reparametrizable, so we can apply $\nabla_{\theta} E_{a \sim \pi_{\theta}}[Q_{\phi}(x, a)] = E_{\Sigma \sim r}[\nabla_{\theta} Q_{\phi}(x, \psi(x, \theta, \epsilon))]$

$$\begin{aligned} \nabla_{\theta} E_{a \sim \pi_{\theta}}[Q_{\phi}(x, a)] &= E_{\Sigma \sim r}[\nabla_{\theta} Q_{\phi}(x, \psi(x, \theta, \epsilon))] \quad a = \psi(x, \theta, \epsilon) \\ &= E_{\Sigma}[\nabla_a Q_{\phi}(x, a)]|_{a=\psi(x, \theta, \epsilon)} \cdot \nabla_{\theta} \psi(x, \theta, \epsilon) \end{aligned}$$

and thus we can compute gradients in this case.

→ SVG (stochastic value gradients) algorithm uses this.

→ SVG often does not explore enough, and lands in local optima.

Maximum entropy reinforcement learning (MERL): Encourages exploration, by letting policies exhibit some uncertainty.

For this we redefine the policy value function to be

$$J_{\lambda}(\theta) = J(\theta) + \lambda H(\pi_{\theta}) = E_{(x,a) \sim \pi_{\theta}}[r(x,a) + \lambda H(\pi(\cdot|x))]$$

λ : Temperature parameter.

The preference for entropy in the actor distribution encourages exploration.

→ We can suitably define regularized (action)-value functions, which we then call "soft" value functions.

→ SAC algorithm: further uses reparametrization gradients.

Language Models as Agents: In some applications (e.g. chatbots) it is difficult to quantify the reward associated with one or a sequence of actions. To align an agent's behavior with human expectations, there are several methods to present feedback to the agent: (1) numerical score. In practice too subjective to give reliable results. (2) Comparison-based feedback. Contains less information, thus takes longer to learn complex behavior. Easy to collect, since humans can easily give comparison feedback.

In language models the agent is given a prompt, and returns a (stochastic) response based on previous responses. The generation of a response can be understood as a policy $\pi_{\theta}(y_{t+1}|x_1, y_1:n)$, which generates the next token (letter/word/...) of the response, depending on the prompt, and all previous tokens.

The training of such models usually consists of 2 main steps.

(1) supervised fine-tuning, training the model on example data.

(2) post-training, using preference feedback

For post-training we can define rewards, based on comparisons

$$y_A \succ y_B | x \quad \text{if } r(y_A|x) > r(y_B|x).$$

A popular choice for modelling preferences is the Bradley-Terry model:

$$\hat{P}(y_A \succ y_B | x, r) = \frac{\exp(r(y_A|x))}{\exp(r(y_A|x)) + \exp(r(y_B|x))} \quad \text{for some } r(y|x)$$

with loss

$$\text{loss}(\hat{P}) = \sum_{y_A, y_B \in \mathcal{Y}} \mu(1) \log(\hat{P}(y_A \succ y_B | x, \hat{r})) + \mu(2) \log(\hat{P}(y_B \succ y_A | x, \hat{r}))$$

In the case of ChatGPT, we first train on a massive dataset, then finetune over high quality data, then learn rewards from comparison feedback, then train with actor critic on learned reward

→ Add penalty term to keep our policy close to the model which was used to generate comparison feedback to train r .

$$\text{Thus use } r(y) = \hat{P}(y) - \lambda \text{KL}(\pi_{\theta_{\text{oppo}}}(y) \| \pi_{\theta_{\text{base}}}(y))$$

→ policy from pre-training/fine tuning.

Model-based Reinforcement Learning: Learning a model can help to dramatically reduce the sample complexity, compared to model-free techniques. For this we will use function approximation to approximate the dynamics model f & p and rewards r , which model our environment. Further, we apply Bayesian Inference, since often understanding the uncertainty in our model is crucial for planning.

The general approach in Model-based RL consists of (1) start with some T_{init} , and initial data D . Then (2) for several iterations (2.1) roll out π to collect observations, (2.2) learn model of f and r from data D , (2.3) plan new policy π based on estimated f, r .

Main challenges: perform planning, learn f and r accurately and efficiently, effectively trade exploration and exploitation.

Planning: In the following we will focus on the setting of continuous state and action spaces, fully observable state spaces, no constraints, and non-linear dynamics model.

Planning with known deterministic dynamics model: For our dynamics model $x_{t+1} = f(x_t, a_t)$, our objective becomes

$$\max_{a_{t:t+\infty}} \sum_{t=0}^{\infty} \gamma^t r(x_t, a_t) \quad \text{s.t. } x_{t+1} = f(x_t, a_t) \quad \rightarrow \begin{array}{l} \text{infinite horizon} \\ \text{discounted return} \end{array}$$

But we cannot explicitly optimize over an infinite horizon.
→ receding horizon control (RHC) / model predictive control (MPC):

To address the above problem, we can iteratively plan over finite horizons i.e. in each round plan over finite horizon H , and then carry out the first action, and after replan.

In particular: (1) at each iteration t , observe X

$$(2) \text{ Optimize performance over: } \max_{a_{t:t+H-1}} \sum_{i=t}^{t+H-1} \gamma^{i-t} r_i(x_i, a_i) \quad \text{s.t. } x_{i+1} = f(x_i, a_i)$$

(3) carry out a_t , then replan.

We observe that X_i can be interpreted as a deterministic function $X_i(a_{t:i-1})$ dependent on all previous actions (from time t). Thus, for (2), we

$$\text{maximize } J_H(a_{t:t+H-1}) = \sum_{i=t}^{t+H-1} \gamma^{i-t} r(X_i(a_{t:i-1}), a_i) \quad \rightarrow \text{in general non-convex.}$$

→ In continuous action space, and differentiable f and r , we can analytically compute gradients (by backpropagating through time) → using the chain rule

→ difficult for large H : local minima, vanishing/exploding gradients.

Random Shooting Methods describe a sampling approach towards global optimization of $J_H(a_{t:t+H-1})$. First generate m sets of random samples $a_{t:t+H-1}^{(i)}$ (e.g. from Gaussian, cross entropy method,...). Then pick the sequence that optimizes J_H . gradually adapts sampling dist. from their rewards.

In finite horizon planning with sparse rewards, H might not go far enough, to get a sense of the optimal reward. To address this, we try estimating also the tail of the discounted return with the value function:

$$J_H(a_{t:t+H-1}) = \underbrace{\sum_{i=t}^{t+H-1} \gamma^{i-t} r(X_i(a_{t:i-1}), a_i)}_{\text{long term}} + \underbrace{\gamma^H V(x_{t+H})}_{\text{short term}}$$

We can use methods like TD-1/Q-Learning to approximate V .

→ Setting $H=1$ coincides with the greedy policy, and we arrive in the model-free setting.

Planning with stochastic dynamics model: trajectory sampling is a generalization of MPC over stochastic models:

- (1) at each iteration t , observe state x_t ,
- (2) Optimize the expected performance over the horizon H :

$$\max_{a_{t:t+H-1}} E_{X_{t+1:t+H}} \left[\sum_{i=t}^{t+H-1} \gamma^{i-t} \cdot r_i + \gamma^H V(x_{t+H}) | a_{t:t+H-1} \right]$$

(3) carry out action a_t , then replan.

Explicitly computing this integral is not feasible. Thus it is common to approximate the expectation using Monte Carlo (trajectory) sampling. But notice that the sampled trajectory depends on the picked actions. To resolve this issue, we can use the reparameterization trick.

We say a dynamics model f is reparametrizable iff: $x_{t+1} = f(x_t, a_t, \varepsilon_t)$ where ε_t is independent of x and a .

Then x_i is determined by $a_{t:i-1}$ and $\varepsilon_{t:i-1}$ via

$$x_i = x_i(a_{t:i-1}, \varepsilon_{t:i-1}) := f(f(\dots(f(x_t, a_t, \varepsilon_t), a_{t+1}, \varepsilon_{t+1}), \dots), a_{i-1}, \varepsilon_{i-1})$$

This allows us to obtain unbiased samples of J_H using Monte Carlo sampling:

$$J_H(a_{t:t+H-1}) \approx \frac{1}{m} \sum_{k=1}^m \sum_{i=1}^{t+H-1} \left(\gamma^{i-1} \cdot r(x_i(\varepsilon_{t:i-1}, a_{t:i-1}), a_i) + \gamma^H V(x_{t+H}) \right)$$

$\downarrow i.i.d.$

→ optimize this via analytic gradient, or shooting methods.

Planning with parametrized policies: When the time horizon is large, or we encounter similar states many times, it can be beneficial to parameterize the policy π directly: $a_t = \pi_\theta(x_t)$

This policy can be trained offline, and evaluated cheaply online. → "open loop control"
Analogously to Q-Learning/DDPG/SVG, the objective becomes

$$J(\theta) = E_{x_0 \sim \mu} \left[\sum_{i=0}^{H-1} \gamma^i r_i + \gamma^H Q_\theta(x_H, \pi_\theta(x_H)) | \theta \right] \quad (\mu(x) > 0 \forall x)$$

→ $H=0$: identical to DDPG objective.

→ Can use reparameterization to allow stochastic policies.

Note: Instead of using Monte Carlo approximation to evaluate a policy, there are more refined ways to approximate $J(\theta)$ (e.g. Moment matching, Variational Inference)

Learn the Dynamics Model: In RL the dynamics model f and rewards r are (of course!) not known.

Key insight: observed transitions and rewards are conditionally independent, because of the Markovian structure of the MDP.

This allows us to treat the estimation of f and r as a simple regression (or density estimation in stochastic case) problem. Thus we can estimate f and r off-policy with standard supervised learning techniques from a replay buffer $D = \{(x_i, a_i, r_i, x_{i+1})\}_i$.
So each observation is a labeled data point. input label

→ Note: input depends on our actions taken.

In the following we focus on learning f (learning r is analogous).

In particular we want to learn probabilistic dynamics models: $x_{t+1} \sim f_\theta(x_t, a_t)$.

Example: Conditional Gaussian dynamics $x_{t+1} \sim \mathcal{W}(\mu_\theta(x_t, a_t), \Sigma_\theta(x_t, a_t))$

represent Σ_θ via lower triangular matrix: $\Sigma_\theta = C_\theta C_\theta^\top$

→ only needs $n(n+1)/2$ param; guarantees automatically p.s.d.

→ allows reparameterization: $x_{t+1} = \mu_\theta(x_t, a_t) + C_\theta(x_t, a_t) \cdot \varepsilon$, $\varepsilon \sim \mathcal{W}(0, I)$.



↳ (learning the dynamics model)

- A first approach to estimating f could be to use the MAP estimate (given a prior and likelihood).
 However, using point estimates leads to a key pitfall of Model-based RL.
 → often performs very poorly.
 → Errors in model estimate compound over multiple time steps ($H \geq 1$).
 These errors are exploited by planning algorithm → overfitting.
 → To solve this pitfall, we need to capture uncertainty in our estimated model. (→ epistemic and aleatoric uncertainty).

Bayesian Learning of the dynamics model: captures the uncertainty about our model by modelling a distribution over f , e.g. model F as a GP or BNN, given a prior belief about f (or r).
 → This allows us to use all (approximate) inference techniques which we learnt earlier. (e.g. variational inference, dropout ensembles, ...)

Let F be chosen from some prior distribution $p(f)$, and we obtain the posterior distribution $P(F|D)$ for $X_{t+1} \sim P(X_t, a_t)$.
 Note that Epistemic: uncertainty in $P(F|D)$ → uncertainty about the model
 Aleatoric: uncertainty in $P(X_{t+1}|f, X_t, a_t)$ → uncertainty about MDP transitions (noise!).
 Recall that epistemic uncertainty corresponds to a distribution over MDPs F . Whereas aleatoric uncertainty corresponds to randomness in transitions within the MDP F .
 → disregard epistemic uncertainty, once a MDP was chosen for planning. → then only concentrate on aleatoric uncertainty.

Thus we receive the following estimate of our reward

$$J_H(a_{t:t+H-1}) \approx \frac{1}{m} \sum_{k=1}^m \sum_{i=t}^{t+H-1} r_i^{i-t} r_i(X_i(a_{t:i-1}, \varepsilon_{t:i-1}^{(k)}, f^{(k)}), a_i) + \gamma^H V(X_{t+H})$$

where $f^{(i)} \sim P(F|D)$ and

$$X_i := X_i(a_{t:i-1}, \varepsilon_{t:i-1}, f) := f(f(\dots(f(x_t, a_t, \varepsilon_t), a_{t+1}, \varepsilon_{t+1}), \dots), a_{i-1}, \varepsilon_{i-1})$$

Using Monte Carlo sampling of $f^{(k)}$! → We are essentially using Monte Carlo Sampling over an "ensemble" of MDP's.

This leads us to a Greedy exploitation algorithm for model-based RL. It greedily maximizes the expected reward w.r.t. the transition model, taking into account epistemic uncertainty.
 Particularly:
 (1) Initialize (possibly empty) data D , and prior $P(F) = p(f|D)$
 (2) For several episodes do the following:
 (2.1) plan new policy π (approx.) maximizing: $\max_{\pi} E_{f \sim P(F|D)} [J_H(\pi; f)]$
 (2.2) roll out π to collect more data in D .
 (2.3) update posterior $P(F|D)$.

PETS (probabilistic ensembles with trajectory sampling) is the above algorithm in the context of NN. It uses an ensemble of NN, each predicting cond. Gaussian transition distributions, and Trajectory sampling to evaluate performance, and MPC for planning.
 → Exploration only happens due to uncertainty in the model.

→ PILCO is algorithm in context of Gaussian Process models. (Moment matching instead of MC sampling)

Exploration in learning dynamics model: A key difference between RL and classical supervised learning is that the chosen actions affect the data which the model learns from.
 → Exploration-exploitation Dilemma.

This problem can be resolved using: added exploration noise, Thompson Sampling, Optimistic exploration.

e.g.
 ↗ Gaussian noise
 "dithering"

For learning the dynamics model

Thompson Sampling: Recall that TS, instead of picking the best action across all realizations of f , it picks the action which performs best for a single realization of f . The epistemic uncertainty in the realizations of f leads to variance in the picked actions, which incentivizes exploration more.

Particularly: (1) Initialize (possibly empty) data D , and prior $p(f) = p(f|D)$.

(2) For several episodes do the following:

(2.1) sample model $\hat{f} \sim p(\cdot|D)$

(2.2) plan policy π , to (approx.) maximize: $\max_{\pi} J_H(\pi; f)$.

(2.3) roll out π , collect more data

(2.4) update posterior $p(f|D)$.

(learning dynamics model)

Bayesian Optimization.

Optimistic Exploration: We have seen in BO that optimism is a central pillar for exploration in (tabular) RL.

Let us consider a set $M(D)$ of plausible models, given data D .
(e.g. $M(D) = \{f_i | f_i(x,a) \in M_i(x,a|D) \pm \beta_0 i(x,a|D), \forall x,a\}$ for cond. Gaussians)

→ This leads us to Optimistic Exploration algorithm, which picks the best action w.r.t. the most optimistic model in $M(D)$.

Particularly: (1) Initialize (possibly empty) data D , and prior $p(f) = P(f|D)$.

(2) for several episodes do the following:

(2.1) plan policy π , to (approx.) maximize: $\max_{\pi} \max_{f \in M(D)} J_H(\pi; f)$

(2.2) roll out π , collect more data

(2.3) update posterior $p(f|D)$.

→ In general, the joint maximization over π and f is very hard.

This leads us to H-UCRL (hallucinated upper confidence reinforcement learning).

To solve the above problem, it interprets the above problem as maximizing the expected reward in an augmented (optimistic) MDP with known dynamics model \hat{f} , with larger action space and additional decision variables ($n(\cdot)$).

Particularly, it considers the optimization problem

$$\pi_{t+1} = \arg \max_{\pi} \max_{n(\cdot) \in [-1, 1]^p} J_H(\pi; \hat{f}_t)$$

with "optimistic dynamics": $\hat{f}_{t,i}(x,a) = M_{t,i}(x,a) + \beta_{t,i} \cdot n_i(x,a) \cdot \sigma_{t,i}(x,a)$

→ n_i controls the variance of an action. Thus, the agent can freely choose the state, which it lands in, taking action a .

→ For this known MDP we can use our developed toolbox.

→ Optimistic exploration works best for hard exploration (e.g. large penalties for certain actions, or sparse rewards).

→ Even for small penalties optimistic explorations learns good policies faster.

Constrained (safe) Exploration: In high-stakes applications, we need to guarantee safety (i.e. avoid unsafe states). Thus exploration can be dangerous.

For this, we denote a set X_{unsafe} of unsafe states, to pessimistically forecast plausible consequences of our actions.

The idea of constrained MDPs is to be optimistic for the rewards, and pessimistic for the costs. But first, optimize

$$\max_{\pi} J(f^*, \pi) = \max_{\pi} E_{x_0 \sim \mu, x_1, \dots \mid \pi, f} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] \quad \begin{matrix} M: \text{distribution over} \\ \text{initial state.} \end{matrix}$$

$$\text{subject to } J_C(f^*, \pi) = E_{x_0 \sim \mu, x_1, \dots \mid \pi, f} \left[\sum_{t=0}^{\infty} \gamma^t \cdot C(x_t) \right] \leq S$$

→ bounds probability of visiting an unsafe state with π .

↓ for some cost function $C(x)$, e.g. $C(x) = \mathbb{I}\{x_{\text{unsafe}}\}$.

→ augmented Lagrangian method is used to solve this.

↓ (Constrained exploration)

→ But note that it is not ensured that constraints are not violated during optimal policy search. Thus, this method is more applicable in an episodic setting. (simulated environment)

Thus we arrive at the strategy to be optimistic wrt. future rewards, and pessimistic wrt. future constraint violations. Solve:

$$\max_{\pi} \max_{\tilde{F} \in M(D)} J(\tilde{F}, \pi) \text{ subject to } \max_{\tilde{F} \in M(D)} J_C(\tilde{F}, \pi) \leq \delta$$

→ We call this the optimistic/pessimistic CMDP

→ LAMBDA: solve the above CMDP using augmented Lagrangian method

- feasible solutions on all tasks
 - better sample efficiency than model-free baselines
 - improvements in safety during training.
- $\left. \begin{array}{l} \text{experimental results} \\ \text{of LAMBDA} \end{array} \right\}$

Safety Filters: To address unsafe exploration we can (slightly) adjust a potentially unsafe policy, to obtain a policy $\hat{\pi}$, which avoids entering unsafe states with high probability. For this we can use the H-UCRL reparametrization to verify the feasibility of any π . Optimize

$$\max_{n(t) \in \mathbb{R}_{+}^{|A|}} J_C(\tilde{F}, [\pi, n]) \leq \delta \quad \text{st. } \tilde{F}(x, a) = M_{t-1}(x, a) + \beta_{t-1} \cdot \sum_{a'} n(a') \cdot R(x, a).$$

and then $n^* = \arg \max_n J_C(\tilde{F}, [\pi, n])$

⇒ Pessimistic value $C_{\pi}^{(P)}(x) = J_C(\tilde{F}, [\pi, n^*] \mid x_0=x)$

Then we can find the "maximally safe" policy via robust RL

$$\pi_{\text{safe}} = \arg \min_{\pi} \max_n J_C(\tilde{F}, [\pi, n])$$

Given an arbitrary policy π , we can apply the safety filter:

$$\Rightarrow \hat{\pi}(x) = \arg \min_a \| \pi(x) - a \|$$

subject to $C_{\pi_{\text{safe}}}^{(P)}(x') \stackrel{\text{next state}}{\leq} \delta$.