

Probabilistic Artificial Intelligence

Cristian Perez Jensen

February 7, 2025

Note that these are not the official lecture notes of the course, but only notes written by a student of the course. As such, there might be mistakes. The source code can be found at github.com/cristianpjensen/eth-cs-notes. If you find a mistake, please create an issue or open a pull request.

Contents

1	Probability review	1
1.1	Random variables	1
1.2	Multivariate Gaussians	2
1.3	Kalman filters	3
1.4	Entropy	4
2	Bayesian linear regression	6
3	Gaussian processes	7
3.1	Learning and inference	8
3.2	Kernel functions	8
3.3	Model selection	9
3.4	Efficiency	10
4	Variational inference	12
4.1	Training	13
4.2	Inference	15
5	Markov chain Monte Carlo	16
5.1	Markov chains	16
5.2	Sampling	18
5.3	Proposal distributions	18
6	Bayesian neural networks	21
6.1	Variational inference	22
6.2	Markov chain Monte Carlo	22
6.3	Monte Carlo dropout	23
6.4	Probabilistic ensembles	23
6.5	Calibration	23
7	Active learning	25
7.1	Sampling strategies	25
8	Bayesian optimization	28
8.1	Acquisition functions	28
9	Markov decision processes	30
9.1	Bellman expectation equation	30
9.2	Policy iteration	31
9.3	Value iteration	31
10	Reinforcement learning	33
10.1	Model-based	33
10.2	Model-free	34
10.3	Model-free deep RL	35
10.4	Model-based deep RL	40

List of symbols

\doteq	Equality by definition
\approx	Approximate equality
\propto	Proportional to
\mathbb{N}	Set of natural numbers
\mathbb{R}	Set of real numbers
$i : j$	Set of natural numbers between i and j . I.e., $\{i, i+1, \dots, j\}$
$f : A \rightarrow B$	Function f that maps elements of set A to elements of set B
$\mathbb{1}\{\text{predicate}\}$	Indicator function (1 if predicate is true, otherwise 0)
$\mathbf{v} \in \mathbb{R}^n$	n -dimensional vector
$\mathbf{M} \in \mathbb{R}^{m \times n}$	$m \times n$ matrix
$\mathbf{T} \in \mathbb{R}^{d_1 \times \dots \times d_n}$	Tensor
\mathbf{M}^\top	Transpose of matrix \mathbf{M}
\mathbf{M}^{-1}	Inverse of matrix \mathbf{M}
$\det(\mathbf{M})$	Determinant of \mathbf{M}
$\frac{d}{dx}f(x)$	Ordinary derivative of $f(x)$ w.r.t. x at point $x \in \mathbb{R}$
$\frac{\partial}{\partial \mathbf{x}}f(\mathbf{x})$	Partial derivative of $f(\mathbf{x})$ w.r.t. \mathbf{x} at point $\mathbf{x} \in \mathbb{R}^n$
$\nabla_{\mathbf{x}}f(\mathbf{x}) \in \mathbb{R}^n$	Gradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $\mathbf{x} \in \mathbb{R}^n$
$\nabla_{\mathbf{x}}^2f(\mathbf{x}) \in \mathbb{R}^{n \times n}$	Hessian of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $\mathbf{x} \in \mathbb{R}^n$
$\boldsymbol{\theta} \in \Theta$	Parametrization of a model, where Θ is a compact subset of \mathbb{R}^K
\mathcal{X}	Input space
\mathcal{Y}	Output space
$\mathcal{D} \subseteq \mathcal{X} \times \mathcal{Y}$	Labeled training data

1 Probability review

Probability is formalized by a probability space (Ω, \mathcal{F}, P) , where Ω is a set of atomic events, $\mathcal{F} \subseteq 2^\Omega$ is the set of non-atomic events, and $P : \mathcal{F} \rightarrow [0, 1]$ is the probability measure that assigns probabilities to events.

The following axioms hold:

$$P(\Omega) = 1 \quad (\text{Normalization})$$

$$P(A) \geq 0 \quad \forall A \in \mathcal{F} \quad (\text{Non-negativity})$$

$$A_1, \dots, A_n \in \mathcal{F} \wedge \bigcap_{i=1}^n A_i = \emptyset \implies P\left(\bigcup_{i=1}^n A_i\right) = \sum_{i=1}^n P(A_i) \quad (\sigma\text{-additivity}).$$

1.1 Random variables

Events are cumbersome to work with, so we can define random variables $X : \Omega \rightarrow D$ for some set D . Then, we can give a probability to X assuming state x ,

$$P(X = x) = P(\{\omega \in \Omega : X(\omega) = x\}).$$

Instead of random variables X , we can also define random vectors $\mathbf{X} = [X_1(\omega), \dots, X_n(\omega)]$. Then, we can specify the joint distribution $P(X_1 = x_1, \dots, X_n = x_n) = P(\mathbf{X} = \mathbf{x})$ succinctly.

For random variables, we have the following rules,

- *Product rule*,

$$P(X_{1:n}) = P(X_1) \prod_{i=2}^n P(X_i \mid X_{1:i-1});$$

- *Sum rule*,

$$P(X_{1:i-1}, X_{i+1:n}) = \sum_{x_i} P(X_{1:i-1}, X_i = x_i, X_{i+1:n});$$

- *Bayes rule*, where we compute the *posterior* $P(X \mid Y)$ from the *likelihood* $P(Y \mid X)$, *prior* $P(X)$, and *marginal* $P(Y)$,

$$P(X \mid Y) = \frac{P(Y \mid X)P(X)}{P(Y)};$$

- A random variable X is *independent* from Y if the following holds for all values,

$$P_{X_1 \dots X_n}(x_1, \dots, x_n) = P_{X_1}(x_1) \cdots P_{X_n}(x_n);$$

- Random variables X and Y are *conditionally independent* given Z if the following holds for all x, y, z ,

$$P_{XY|Z}(x, y \mid z) = P_{X|Z}(x \mid z)P_{Y|Z}(y \mid z).$$

1.2 Multivariate Gaussians

Suppose we have n binary variables, then we need $2^n - 1$ parameters.^{1,2} Also, if we want to compute the joint distribution over all n variables,³ we would have to sum up 2^{n-1} terms according to the sum rule. In conclusion, binary random variables scale poorly. Furthermore, we would need a lot of data to estimate the distribution.

The solution to these problems are multivariate Gaussians,

$$\mathcal{N}(x; \mu, \Sigma) = \frac{1}{2\pi\sqrt{\det(\Sigma)}} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right),$$

where $\mu \in \mathbb{R}^n$ and $\Sigma \in \mathbb{S}_{++}^n$.⁴ Thus, the joint distribution over n Gaussian variables requires only $n^2 + n$ parameters.

Let $X \sim \mathcal{N}(\mu, \Sigma)$ be a d -dimensional Gaussian random vector, then the following properties hold,

- Let A be an index set, then the marginal distribution of variables indexed by A is the following,

$$X_A \sim \mathcal{N}(\mu_A, \Sigma_{AA}).$$

Thus, it is simply a look-up to get a subset marginal distribution;

- Let A and B be index sets, then the marginal distribution of variables indexed by A , conditioned on B , is the following,

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

where

$$\begin{aligned} \mu_{A|B} &= \mu_A + \Sigma_{AB}\Sigma_{BB}^{-1}(x_B - \mu_B) \\ \Sigma_{A|B} &= \Sigma_{AA} - \Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA}. \end{aligned}$$

Notice that $\Sigma_{AB}\Sigma_{BB}^{-1}$ removes the interdependencies of B and adds the dependencies of B with A . Further notice that the dependency of A and B added scales linearly with the mean difference $x_B - \mu_B$.

Further, notice that $\Sigma_{A|B}$ only depends on which random variables are observed, not what values those random variables are, because it does not depend on x_B ;

- Let $M \in \mathbb{R}^{m \times d}$ be a matrix, then $Y = MX$ is also a Gaussian,

$$Y \sim \mathcal{N}(M\mu, M\Sigma M^\top).$$

Notice that m is not necessarily equal to d , so we can transform a d -dimensional random vector to any dimensionality m ;

- Let X' be another d -dimensional Gaussian, then $Y = X + X'$ is also a Gaussian,

$$Y \sim \mathcal{N}(\mu + \mu', \Sigma + \Sigma').$$

¹ -1 , because we do not need to specify the last parameter, since it will be whatever is remaining of the total probability.

² In other words, the parametrization of the distribution grows exponentially.

³ I.e., do inference.

⁴ Σ is an $n \times n$ positive semi-definite matrix.



Figure 1.1. Bivariate Gaussian distribution with $\Sigma = I$.



Figure 1.2. Bivariate Gaussian distribution with

$$\Sigma = \begin{bmatrix} 1 & 1/2 \\ 1/2 & 1 \end{bmatrix}.$$

If x_1 increases, the probability of a higher x_2 increases as well.

1.3 Kalman filters

Definition 1.1 (Kalman filter). A Kalman filter is specified by a Gaussian prior over the states,

$$X_0 \sim \mathcal{N}(\mu, \Sigma),$$

and a conditional linear Gaussian *motion model* and *sensor model*,

$$\begin{aligned} X_{t+1} &\doteq FX_t + \epsilon_t & F &\in \mathbb{R}^{d \times d}, & \epsilon_t &\sim \mathcal{N}(\mathbf{0}, \Sigma_x) \\ Y_t &\doteq HX_t + \eta_t & H &\in \mathbb{R}^{m \times d}, & \eta_t &\sim \mathcal{N}(\mathbf{0}, \Sigma_y), \end{aligned}$$

respectively.



Figure 1.3. Directed graphical model of a Kalman filter with hidden states X_t and observables Y_t .

From the directed graphical model in Figure 1.3, we can observe the following conditional independences,

$$\begin{aligned} X_{t+1} &\perp\!\!\!\perp X_{1:t-1}, Y_{1:t-1} \mid X_t \\ Y_t &\perp\!\!\!\perp X_{1:t-1} \mid X_t \\ Y_t &\perp\!\!\!\perp Y_{1:t-1} \mid X_{t-1}. \end{aligned}$$

These lead to the following factorization of the joint distribution,

$$\begin{aligned} p(\mathbf{x}_{1:t}, \mathbf{y}_{1:t}) &= \prod_{i=1}^t p(x_i \mid \mathbf{x}_{1:i-1}) p(y_i \mid \mathbf{x}_{1:t}, \mathbf{y}_{1:i-1}) \\ &= p(x_1) p(y_1 \mid x_1) \prod_{i=2}^t p(x_i \mid x_{i-1}) p(y_i \mid x_i). \end{aligned}$$

We now want to do *Bayesian filtering* on a Kalman filter, which involves keeping track of an agent's state using noisy observations Y . It is described by the recursive scheme in Figure 1.4.



Figure 1.4. The recursive scheme of Bayesian filtering.

We can do the update by the following,

$$\begin{aligned} p(x_t | \mathbf{y}_{1:t}) &= \frac{1}{Z} p(x_t | \mathbf{y}_{1:t-1}) p(y_t | x_t, \mathbf{y}_{1:t-1}) \\ &= \frac{1}{Z} p(x_t | \mathbf{y}_{1:t-1}) p(y_t | x_t). \end{aligned}$$

Furthermore, we can do the prediction by the following,

$$\begin{aligned} p(x_{t+1} | \mathbf{y}_{1:t}) &= \int p(x_{t+1}, x_t | \mathbf{y}_{1:t}) dx_t \\ &= \int p(x_{t+1} | x_t, \mathbf{y}_{1:t}) p(x_t | \mathbf{y}_{1:t}) dx_t \\ &= \int p(x_{t+1} | x_t) p(x_t | \mathbf{y}_{1:t}) dx_t. \end{aligned}$$

In general, these distributions are very complicated, but for Gaussians (as in the Kalman filter), they can be expressed in closed form. The general formula for the *Kalman update* is as follows, given the prior belief $X_t | \mathbf{y}_{1:t} \sim \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$,

$$\begin{aligned} X_{t+1} | \mathbf{y}_{1:t+1} &\sim \mathcal{N}(\boldsymbol{\mu}_{t+1}, \boldsymbol{\Sigma}_{t+1}) \\ \boldsymbol{\mu}_{t+1} &\doteq \mathbf{F}\boldsymbol{\mu}_t + \mathbf{K}_{t+1}(\mathbf{y}_{t+1} - \mathbf{H}\mathbf{F}\boldsymbol{\mu}_t) \\ \boldsymbol{\Sigma}_{t+1} &\doteq (\mathbf{I} - \mathbf{K}_{t+1}\mathbf{H})(\mathbf{F}\boldsymbol{\Sigma}_t\mathbf{F}^\top + \boldsymbol{\Sigma}_x), \end{aligned}$$

where \mathbf{K}_{t+1} is the *Kalman gain*,

$$\mathbf{K}_{t+1} \doteq (\mathbf{F}\boldsymbol{\Sigma}_t\mathbf{F}^\top + \boldsymbol{\Sigma}_x)\mathbf{H}^\top (\mathbf{H}(\mathbf{F}\boldsymbol{\Sigma}_t\mathbf{F}^\top + \boldsymbol{\Sigma}_x)\mathbf{H}^\top + \boldsymbol{\Sigma}_y)^{-1}.$$

The term $(\mathbf{y}_{t+1} - \mathbf{H}\mathbf{F}\boldsymbol{\mu}_t)$ measures the error in the predicted observation and the Kalman gain \mathbf{K}_{t+1} measures the relevance of the new observation compared to the prediction.

1.4 Entropy

Definition 1.2 (Entropy). *Entropy* measures the expected surprisal of a distribution p ,

$$H[p] \doteq \mathbb{E}_{x \sim p}[-\log p(x)].$$

$-\log p(x)$ is also called the surprisal value of x , because it decreases as the probability grows, and is exactly 0 if the probability is 1.

Definition 1.3 (Kullback-Leibler divergence). *Kullback-Leibler divergence* (KL divergence) is a common metric that measures dissimilarity between two distributions p and q ,

$$KL(p||q) \doteq \mathbb{E}_{\theta \sim p} \left[\log \frac{p(\theta)}{q(\theta)} \right].$$

Intuitively, $KL(p||q)$ measures the information loss when approximating p with q .

Definition 1.4 (Mutual information). Given random variables X and Y , the mutual information $I(X; Y)$ quantifies how much observing Y reduces uncertainty about X , as measured by its entropy, in expectation over Y .

$$I(X; Y) \doteq H[X] - H[X | Y],$$

where $H[X]$ and $H[X | Y]$ quantify the uncertainty about X before and after observing Y .

Properties (Mutual information). Mutual information is symmetric,

$$I(X; Y) = I(Y; X).$$

Mutual information is positive (*information never hurts*),

$$I(X; Y) \geq 0.$$

Example 1.5 (Mutual information of noisy Gaussian observations). Let

$$\begin{aligned} X &\sim \mathcal{N}(\mu, \Sigma) \\ Y &= X + \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \sigma_n^2 \mathbf{I}). \end{aligned}$$

Then,

$$\begin{aligned} I(X; Y) &= H[Y] - H[Y | X] \\ &= H[Y] - H[\epsilon] \\ &= \frac{1}{2} \log(2\pi e)^d \det(\Sigma + \sigma_n^2 \mathbf{I}) - \frac{1}{2} \log(2\pi e)^d (\sigma_n^2 \mathbf{I}) \\ &= \frac{1}{2} \log \det(\mathbf{I} + \sigma_n^{-2} \Sigma). \end{aligned}$$

Definition 1.6 (Conditional mutual information).

$$\begin{aligned} I(X; Y | Z) &\doteq H[X | Z] - H[X | Y, Z] \\ &= I(X; Y, Z) - I(X; Z). \end{aligned}$$

2 Bayesian linear regression

Bayesian linear regression (BLR) is a model that is able to provide an uncertainty measure about its predictions due to a lack of data.⁵ It does so by outputting a probability distribution over possible outputs y^* given input x^* . The variance of this distribution measures the uncertainty of the model at this point and the mode corresponds to its best estimate. It does this by not considering a single set of weights, but rather all possible weights, assigning each a probability via Bayes rule. Recall the posterior,

$$p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}) \propto p(\mathbf{y} \mid \mathbf{w}, \mathbf{X}) \cdot p(\mathbf{w}, \mathbf{X}),$$

where we can compute the maximum a posteriori (MAP) estimate and use that as our weights, multiplying it with the input x^* . However, by considering every plausible function and giving it a probability according to how well it models the data, we can compute the uncertainty of the model. Instead of using only the mode of the posterior, we will use the full posterior of \mathbf{w} by taking the integral over all of them. This allows us to assign a probability distribution to any point (x^*, y^*) ,

$$\begin{aligned} p(y^* \mid x^*, \mathbf{X}, \mathbf{y}) &= \int p(\mathbf{w}, y^* \mid x^*, \mathbf{X}, \mathbf{y}) d\mathbf{w} && \text{Sum rule} \\ &= \int p(\mathbf{w} \mid x^*, \mathbf{X}, \mathbf{y}) \cdot p(y^* \mid \mathbf{w}, x^*, \mathbf{X}, \mathbf{y}) d\mathbf{w} && \text{Product rule} \\ &= \int p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}) \cdot p(y^* \mid \mathbf{w}, x^*) d\mathbf{w}. \end{aligned}$$

Notice that each model \mathbf{w} 's prediction $p(y^* \mid \mathbf{w}, x^*)$ is weighted by its probability $p(\mathbf{w} \mid \mathbf{X}, \mathbf{y})$. In general, this integral is intractable, but by assuming that the prior and likelihood are independently Gaussian,

$$\begin{aligned} \mathbf{w} &\sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I}) \\ y_i \mid \mathbf{w}, x_i &\sim \mathcal{N}(\mathbf{w}^\top x_i, \sigma_n^2), \end{aligned}$$

it becomes tractable. Firstly, the posterior is given by the following,⁶

$$\begin{aligned} \mathbf{w} \mid \mathbf{X}, \mathbf{y} &\sim \mathcal{N}(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\Sigma}}) \\ \bar{\boldsymbol{\mu}} &= (\mathbf{X}^\top \mathbf{X} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} \\ \bar{\boldsymbol{\Sigma}} &= (\sigma_n^{-2} \mathbf{X}^\top \mathbf{X} + \mathbf{I})^{-1}. \end{aligned}$$

By taking advantage of the Gaussian distribution properties, we can compute the distribution of y^* given a point x^* . Let's say $f^* = \mathbf{w}^\top x^* = x^{*\top} \mathbf{w}$, then

$$f^* \mid x^*, \mathbf{X}, \mathbf{y} \sim \mathcal{N}(x^{*\top} \bar{\boldsymbol{\mu}}, x^{*\top} \bar{\boldsymbol{\Sigma}} x^*).$$

Adding aleatoric noise $y^* = f^* + \epsilon_n$ results in the following,

$$y^* \mid x^*, \mathbf{X}, \mathbf{y} \sim \mathcal{N}(x^{*\top} \bar{\boldsymbol{\mu}}, x^{*\top} \bar{\boldsymbol{\Sigma}} x^* + \sigma_n^2).$$

The epistemic uncertainty $\bar{\boldsymbol{\Sigma}}$ is uncertainty about the model due to a lack of data, while the aleatoric uncertainty σ_n^2 is irreducible noise that is always present in data.

⁵ This is called epistemic uncertainty.

⁶ Notice that $\bar{\boldsymbol{\Sigma}}$ does not depend on \mathbf{y} . Since the covariance matrix measures the uncertainty, this tells us that the uncertainty only depends on where we observed data, not what we observed. Intuitively, this makes a lot of sense.

3 Gaussian processes

BLR can only make linear predictions, because it is linear in the parameters. However, we might want to make non-linear predictions. We could apply BLR on non-linearly transformed data,⁷ but the computational cost would increase with the dimensionality of the feature space.

An alternative way of looking at it is considering inference directly in function space. We use *Gaussian processes* (GP) to describe distributions over functions. They are formally defined as an infinite collection of random variables, of which any finite number have a joint Gaussian distribution. GPs are specified by a mean function $\mu : \mathcal{X} \rightarrow \mathbb{R}$ and a kernel function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, and written as the following,⁸

$$f \sim \mathcal{GP}(\mu, k).$$

Now, we would want a finite Gaussian distribution over a marginal $\{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subseteq \mathcal{X}$ (stored by matrix \mathbf{X}) of the infinite collection of random variables \mathcal{X} . If we assume a normal prior on the weights,

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

then the distribution over $\mathbf{f} = \mathbf{w}^\top \mathbf{X}^\top = \mathbf{X}\mathbf{w}$ becomes the following,

$$\mathbf{f} \sim \mathcal{N}(\mathbf{X}^\top \mathbf{0}, \mathbf{X}^\top \mathbf{I} \mathbf{X}) = \mathcal{N}(\mathbf{0}, \mathbf{X}^\top \mathbf{X}).$$

Notice that the data points enter as inner products, thus we do not necessarily need to let them linearly depend. We could also use any kernel function,⁹

$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, k(\mathbf{X}, \mathbf{X}))$$

with

$$k(\mathbf{X}, \mathbf{X}) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_m) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_m, \mathbf{x}_1) & \cdots & k(\mathbf{x}_m, \mathbf{x}_m) \end{bmatrix}.$$

We can sample function realizations from a Gaussian process, by taking n equidistant points from \mathcal{X} as matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$. Then, we assume $\mu(\mathbf{x}) = 0$ and compute $k(\mathbf{X}, \mathbf{X})$ to make a multivariate Gaussian distribution over $\mathbf{f} \in \mathbb{R}^n$. Lastly, we sample a vector from this multivariate Gaussian distribution and interpolate between the points to form the function realization. Notice that GPs parametrize a probability distribution over functions.

Figure 3.1 shows samples of a prior Gaussian process with the periodic kernel,

$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp\left(-\frac{2}{\ell^2} \sin^2\left(\pi \frac{|\mathbf{x} - \mathbf{x}'|}{p}\right)\right),$$

where σ^2 is the overall variance, ℓ is the lengthscale, and p is the period. It also shows the periodic kernel function w.r.t. $x = 0$. As can be seen,

⁷ This would mean redefining $f(\mathbf{x})$ to

$$f(\mathbf{x}) = \boldsymbol{\phi}(\mathbf{x})^\top \mathbf{w},$$

with e.g.

$$\boldsymbol{\phi}(\mathbf{x}) = [1 \quad x \quad x^2]^\top.$$

⁸ In other words, we sample functions from a GP defined by a mean and kernel function.

⁹ Formally, $k(\mathbf{x}, \mathbf{x}') = \boldsymbol{\phi}(\mathbf{x})^\top \boldsymbol{\phi}(\mathbf{x}')$ for some feature function $\boldsymbol{\phi}$. But, using a kernel function k instead makes it more computationally efficient, because the dimensionality of the Gaussian does not scale with the output dimensionality of the feature function.



Figure 3.1. A priori samples of a Gaussian process with the periodic kernel. The second plot shows the kernel function w.r.t. $x = 0$.

pairs of points (x, x') with high covariance $k(x, x')$ have close values in all sampled function realizations. This makes sense, since k parametrizes the covariance.

3.1 Learning and inference

Adding aleatoric noise, we define the observed data \mathbf{y} to have the following distribution,

$$\mathbf{y} \sim \mathcal{N}(\mathbf{0}, k(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I}).$$

Now, suppose that we observe data \mathbf{y} for datapoints \mathbf{X} , and want to predict the probability distribution of y^* for x^* given the observed data. We can define the following a priori joint distribution,¹⁰

$$\begin{bmatrix} \mathbf{y} \\ f^* \end{bmatrix} \mid \mathbf{x}^*, \mathbf{X} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} k(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I} & k(\mathbf{X}, \mathbf{x}^*) \\ k(\mathbf{x}^*, \mathbf{X}) & k(\mathbf{x}^*, \mathbf{x}^*) \end{bmatrix}\right),$$

where

$$k(\mathbf{X}, \mathbf{x}) = \begin{bmatrix} k(x_1, x) \\ \vdots \\ k(x_m, x) \end{bmatrix}.$$

Then, we can derive the conditional distribution, using the conditional property of multivariate Gaussian distributions, as

$$f^* \mid \mathbf{x}^*, \mathbf{X}, \mathbf{y} \sim \mathcal{N}(\mu^*, k^*),$$

where

$$\begin{aligned} \mu^* &= k(\mathbf{x}^*, \mathbf{X})(k(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y} \\ k^* &= k(\mathbf{x}^*, \mathbf{x}^*) - k(\mathbf{x}^*, \mathbf{X})(k(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I})^{-1} k(\mathbf{X}, \mathbf{x}^*). \end{aligned}$$

Adding aleatoric noise, we get the probability distribution over y^* :

$$y^* \mid \mathbf{x}^*, \mathbf{X}, \mathbf{y} \sim \mathcal{N}(\mu^*, k^* + \sigma_n^2).$$

3.2 Kernel functions

Suppose we have two covariance functions,

$$k_1 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}, \quad k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R},$$

$c > 0$, and f is a polynomial with positive coefficients or the exponential function. Then, the following functions are valid covariance functions,

$$\begin{aligned} k(x, x') &= k_1(x, x') + k_2(x, x') \\ k(x, x') &= k_1(x, x') \cdot k_2(x, x') \\ k(x, x') &= c \cdot k_1(x, x') \\ k(x, x') &= f(k_1(x, x')). \end{aligned}$$

¹⁰ Notice that this is a priori because we have not observed (conditioned on) any data points yet. We are only specifying the joint distribution over \mathbf{y} and f^* .



Figure 3.2. The periodic kernel function hyperparameters used are $\sigma^2 = 1$, $\ell = 1$, and $p = 3$. The second plot shows the covariance w.r.t. $x = 0$.

Definition 3.1 (Stationary and isotropic kernels). A kernel function k is *stationary* if its function only depends on the difference between its arguments, i.e., $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x} - \mathbf{x}')$. It is *isotropic* if it only depends on the ℓ_2 -distance between its arguments, i.e., $k(\mathbf{x}, \mathbf{x}') = k(\|\mathbf{x} - \mathbf{x}'\|_2)$

The following is a list of popular kernels,

- *Linear kernel*,

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}' + \sigma_0^2.$$

Its posterior can be seen in Figure 3.3;

- *Gaussian kernel* (a.k.a. *Squared Exponential* or *RBF*),

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\ell^2}\right).$$

Points that are close together have a high covariance, while points further away have a lower one. This is what makes it smooth, and continuous. Its posterior can be seen in Figure 3.4;

- *Exponential kernel*,

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|}{\ell}\right).$$

Points that are close together have a high covariance. However, because it uses the ℓ_1 -distance, the covariance function is not smooth, but “linear on two sides”. This is why it has so many peaks. Its posterior can be seen in Figure 3.5.



Figure 3.3. Posterior linear kernel. The second plot shows the covariance w.r.t. several points.

3.3 Model selection

As can be seen in Figures 3.2 to 3.5, the hyperparameters matter a lot for whether a GP can model the datapoints correctly. These hyperparameters can be learned by maximizing its predictive performance on the data. (Hyperparameter fitting on the training data does not cause overfitting. This will become clear.)

Suppose we have data $\{(\mathbf{x}_i^*, y_i^*)\}_{i=1}^n$, then we would like to choose hyperparameters, such that the performance on this data is maximized. There are several choices for measuring predictive performance. The most naive option is mean squared error,

$$\text{MSE}_\theta(y^*, \mathbf{x}^*) = (y^* - \mu_\theta^*(\mathbf{x}^*))^2.$$

This metric ignores aleatoric (σ_n^2) and epistemic ($k(\mathbf{x}^*, \mathbf{x}^*)$) uncertainty, thus it will not work well for measuring the performance of GPs. Another option is to just add the variance to the loss, so it also will be minimized,

$$\tilde{\text{MSE}}_\theta(y^*, \mathbf{x}^*) = \text{MSE}_\theta(y^*, \mathbf{x}^*) + k_\theta^*(\mathbf{x}^*, \mathbf{x}^*).$$

The problem with this is that it encourages low epistemic uncertainty which will cause the loss function to favor models that have low epistemic uncertainty without it necessarily being true. Furthermore, it still ignores the aleatoric uncertainty.

The Bayesian perspective provides an alternative approach: maximizing the likelihood of the data,

$$\begin{aligned}\ell\ell_\theta(\mathbf{y}^*, \mathbf{x}^*) &= \mathcal{N}(\mathbf{y}^*; f_\theta^*(\mathbf{x}^*), \sigma_n^2) \\ &= \mathcal{N}(\mathbf{y}^*; \mu_\theta^*, k_\theta^*(\mathbf{x}^*, \mathbf{x}^*) + \sigma_n^2).\end{aligned}$$

We can optimize θ as follows,

$$\begin{aligned}\hat{\theta} &= \operatorname{argmax}_{\theta} p(\mathbf{y}^* | \mathbf{X}^*, \theta) \\ &= \operatorname{argmax}_{\theta} \int p(\mathbf{y}^* | f, \mathbf{X}^*, \theta) p(f | \theta) df \\ &= \operatorname{argmax}_{\theta} \mathcal{N}(\mathbf{y}^*; \mathbf{0}, k(\mathbf{X}^*, \mathbf{X}^*) + \sigma_n^2 \mathbf{I}) \\ &= \operatorname{argmin}_{\theta} -\log \mathcal{N}(\mathbf{y}^*; \mathbf{0}, k(\mathbf{X}^*, \mathbf{X}^*) + \sigma_n^2 \mathbf{I}) \\ &= \operatorname{argmin}_{\theta} \frac{n}{2} \log 2\pi + \frac{1}{2} \log \det(k(\mathbf{X}^*, \mathbf{X}^*) + \sigma_n^2 \mathbf{I}) + \frac{1}{2} \mathbf{y}^{*\top} (k(\mathbf{X}^*, \mathbf{X}^*) + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}^* \\ &= \operatorname{argmin}_{\theta} \underbrace{\frac{1}{2} \log \det(k(\mathbf{X}^*, \mathbf{X}^*) + \sigma_n^2 \mathbf{I})}_{\text{complexity penalty}} + \underbrace{\frac{1}{2} \mathbf{y}^{*\top} (k(\mathbf{X}^*, \mathbf{X}^*) + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}^*}_{\text{"goodness" of fit}}.\end{aligned}$$

As can be seen, the loss function seeks a balance between the “goodness” of the fit and complexity. If we increase the aleatoric uncertainty σ_n^2 , we increase the goodness of the fit, but increase the complexity. This is how it prevents over- and underfitting and is the reason why we do not need a validation dataset.

3.4 Efficiency

The time complexity of computing the posterior is $\Theta(n^3)$ ¹¹ and the space complexity of storing $k(\mathbf{X}, \mathbf{X})$ is $\Theta(n^2)$. The main approaches for accelerating GP posterior computation are exploiting parallelism (GPU),¹² local GP methods, kernel function approximations, and inducing point methods.

Local GP methods. The basic idea is that, for covariance functions that decay with distance,¹³ we only need to condition on close points to \mathbf{x} . *I.e.*, to make a prediction at point \mathbf{x} , we only need to condition on points \mathbf{x}' where $|k(\mathbf{x}, \mathbf{x}')| \geq \tau$ for some threshold τ . The problem with this method is that it is still expensive if there are many close points.

Kernel function approximation. The key idea of approximating kernel functions is that we can construct an m -dimensional feature map with

¹¹ For BLR, this is $O(dn^2)$.

¹² This yields a significant speedup, but does not address the cubic scaling in n .

¹³ Think of stationary kernels such as the Gaussian and exponential kernels.

$m \ll n$ that approximates the true kernel function,

$$k(x, x') \approx \phi(x)^\top \phi(x') \quad \phi(x) \in \mathbb{R}^m.$$

Then, apply BLR. The computational cost becomes $\mathcal{O}(nm^2 + m^3)$ instead of $\mathcal{O}(n^3)$.

An example is Random Fourier Features (RFF) that reduces stationary kernels to their Fourier transform. Then, samples from this m times and defines it as the feature map ϕ . The problem with RFFs is that they approximate the kernel function globally, however this might not be necessary, since we only need accurate representation for the training and test points.

Inducing point methods. The idea behind inducing point methods is that we won't need all data. In areas where there are a lot of data points, we can safely throw some away. This method needs to figure out which points are safe to throw away, *i.e.*, find inducing points \mathcal{U} that we need to keep to approximate well. This can be done by choosing them randomly or we could treat \mathcal{U} as hyperparameters and maximize the marginal likelihood of the data.



Figure 3.4. Posterior Gaussian kernel. The second plot shows the covariance w.r.t. several points.



Figure 3.5. Posterior exponential kernel. The second plot shows the covariance w.r.t. several points.

4 Variational inference

Remark. From now on, \mathbf{X} will be omitted and treated as a constant.

In Bayesian learning, we want to compute the following probability distribution to make predictions given the data,

$$p(y^* | \mathbf{x}^*, \mathbf{y}) = \int p(y^* | \mathbf{x}^*, \theta) p(\theta | \mathbf{y}) d\theta,$$

where we marginalize over all possible models θ . Furthermore, we compute the posterior as follows,

$$p(\theta | \mathbf{y}) = \frac{1}{Z} p(\theta) \prod_{i=1}^n p(y_i | \mathbf{x}_i, \theta).$$

However, in general, these equations are intractable.¹⁴ Gaussian processes solved this problem by assuming that the prior and likelihood are Gaussian.¹⁵ But, in some cases, it is not realistic to assume Gaussian distributions.¹⁶ *Variational inference* solves this problem by approximating the intractable distribution p by a simpler one q that is “as close as possible”,

$$p(\theta | \mathbf{y}) = \frac{1}{Z} p(\theta, \mathbf{y}) \approx q(\theta | \lambda) = q_\lambda(\theta),$$

where λ are called the variational parameters. Thus, we have reduced the problem to optimization, *i.e.*, maximizing the similarity between distributions p and q , where q is part of a variational family \mathcal{Q} that is easy to work with.

¹⁴ The integral is intractable, because distributions are not conjugate in general. The posterior is intractable, because of the normalizer Z .

¹⁵ Because the Gaussian has a conjugate prior.

¹⁶ *E.g.* in logistic regression, the likelihood is modeled by a Bernoulli distribution.

Example 4.1 (Laplace approximation). A simple way of approximating intractable integrals is *Laplace approximation* which is a Gaussian approximation to the posterior. Let's define a function $\psi(\theta) \doteq \log p(\theta | \mathbf{y})$, then the Laplace approximation of ψ can be computed from the second-order Taylor expansion around the posterior mode,

$$\begin{aligned}\psi(\theta) &\approx \psi(\hat{\theta}) + (\theta - \hat{\theta})^\top \nabla \psi(\hat{\theta}) + \frac{1}{2}(\theta - \hat{\theta})^\top \mathbf{H}_\psi(\hat{\theta})(\theta - \hat{\theta}) \\ &= \psi(\hat{\theta}) + \frac{1}{2}(\theta - \hat{\theta})^\top \mathbf{H}_\psi(\hat{\theta})(\theta - \hat{\theta}).\end{aligned}$$

Then,

$$\begin{aligned}p(\theta | \mathbf{y}) &= \exp(\psi(\theta)) \\ &\approx \exp\left(\psi(\hat{\theta}) + \frac{1}{2}(\theta - \hat{\theta})^\top \mathbf{H}_\psi(\hat{\theta})(\theta - \hat{\theta})\right) \\ &= \exp(\psi(\hat{\theta})) \cdot \exp\left(\frac{1}{2}(\theta - \hat{\theta})^\top \mathbf{H}_\psi(\hat{\theta})(\theta - \hat{\theta})\right) \\ &= \frac{1}{Z} \cdot \exp\left(-\frac{1}{2}(\theta - \hat{\theta})^\top \mathbf{H}_\psi^{-1}(\hat{\theta})(\theta - \hat{\theta})\right) \\ &= \mathcal{N}(\theta; \hat{\theta}, \mathbf{\Lambda}^{-1}) \\ &\doteq q(\theta),\end{aligned}$$

where $\hat{\theta} = \operatorname{argmax}_\theta p(\theta | \mathbf{y})$ and $\mathbf{\Lambda} = -\mathbf{H}_\theta \log p(\theta | \mathbf{y})$

Intuitively, the Laplace approximation matches the shape of the true posterior around its mode, but may not represent it accurately elsewhere. Often, this leads to extremely overconfident predictions.

4.1 Training

When training, we want to find the distribution in the variational family \mathcal{Q} that minimizes the KL divergence with p ,¹⁷

¹⁷ When approximating p with q , we can either minimize the reverse or the forward KL divergence. The reverse KL divergence $KL(q||p)$ typically acts more greedily and places most of its mass where p has a lot of mass, while the forward KL divergence $KL(p||q)$ tries to cover most of the probability mass of p . Thus, the forward KL divergence is more desirable, but it requires sample from p , which is intractable (the whole reason we are doing this). Thus, we have to resort to using the reverse KL divergence.

$$\begin{aligned}
q^* &= \operatorname{argmin}_{q_\lambda \in \mathcal{Q}} KL(q_\lambda \| p) \\
&= \operatorname{argmin}_\lambda \int q_\lambda(\theta) \log \frac{q_\lambda(\theta)}{p(\theta | \mathbf{y})} d\theta \\
&= \operatorname{argmin}_\lambda \int q_\lambda(\theta) \log \frac{q_\lambda(\theta)}{\frac{1}{Z} p(\mathbf{y} | \theta) p(\theta)} d\theta \\
&= \operatorname{argmax}_\lambda \int q_\lambda(\theta) (\log p(\mathbf{y} | \theta) + \log p(\theta) - \log Z - \log q_\lambda(\theta)) d\theta \\
&= \operatorname{argmax}_\lambda \int q_\lambda(\theta) \left(\log p(\mathbf{y} | \theta) - \log \frac{q_\lambda(\theta)}{p(\theta)} \right) d\theta \\
&= \operatorname{argmax}_\lambda \int q_\lambda(\theta) \log p(\mathbf{y} | \theta) - \int q_\lambda(\theta) \log \frac{q_\lambda(\theta)}{p(\theta)} d\theta \\
&= \operatorname{argmax}_\lambda \underbrace{\mathbb{E}_{\theta \sim q_\lambda} [\log p(\mathbf{y} | \theta)]}_{\text{expected likelihood}} - \underbrace{KL(q_\lambda \| p_{\text{prior}})}_{\text{"stay close to prior"}}.
\end{aligned}$$

Thus, minimizing $KL(q_\lambda \| p)$ is equivalent to maximizing the expected likelihood, while remaining close to the prior distribution.

To show that minimizing the KL divergence is an adequate method of model selection, we will show that it lower bounds the evidence,

$$\begin{aligned}
\log p(\mathbf{y}) &= \log \int p(\mathbf{y}, \theta) d\theta && \text{Sum rule} \\
&= \log \int q_\lambda(\theta) \frac{p(\mathbf{y}, \theta)}{q_\lambda(\theta)} d\theta \\
&= \log \mathbb{E}_{\theta \sim q_\lambda} \left[\frac{p(\mathbf{y}, \theta)}{q_\lambda(\theta)} \right] \\
&\geq \mathbb{E}_{\theta \sim q_\lambda} \left[\log \frac{p(\mathbf{y}, \theta)}{q_\lambda(\theta)} \right] && \text{Jensen's inequality} \\
&= \mathbb{E}_{\theta \sim q_\lambda} \left[\log \frac{p(\mathbf{y} | \theta) p(\theta)}{q_\lambda(\theta)} \right] \\
&= \mathbb{E}_{\theta \sim q_\lambda} [\log p(\mathbf{y} | \theta)] - KL(q_\lambda \| p_{\text{prior}}).
\end{aligned}$$

Thus, minimizing

$$L(\lambda) = \mathbb{E}_{\theta \sim q_\lambda} [\log p(\mathbf{y} | \theta)] - KL(q_\lambda \| p_{\text{prior}})$$

can safely be used to find an appropriate model for the data, since a lower bound on the evidence gets maximized.

However, there is one problem: we want to compute the gradient of an expectation w.r.t. q , but q depends on λ . Thus, we cannot compute the gradient in its current form. To solve this we use the reparametrization trick.

Definition 4.2 (Reparameterization trick). Suppose we have a random variable $\epsilon \sim \phi$ sampled from a base distribution, and consider $\theta \doteq g(\epsilon, \lambda)$ for some invertible function g . Then, the following holds:

$$\mathbb{E}_{\theta \sim q_\lambda}[f(\theta)] = \mathbb{E}_{\epsilon \sim \phi}[f(g(\epsilon, \lambda))].$$

Thus, after reparametrization, the expectation is w.r.t. to distribution ϕ that does not depend on λ . Thus, we can compute the gradient $\nabla_\lambda L(\lambda)$.

Example 4.3 (Reparametrization trick for Gaussians). Suppose we use a Gaussian variational approximation,

$$q_\lambda(\theta) \doteq \mathcal{N}(\theta; \mu, \Sigma).$$

Then, we can reparametrize $\theta = g(\epsilon, \lambda) = \Sigma^{1/2}\epsilon + \mu$, where $\phi = \mathcal{N}(\mathbf{0}, \mathbf{I})$, because

$$\theta = \Sigma^{1/2}\epsilon + \mu \sim \mathcal{N}(\mu, \Sigma^{1/2}\mathbf{I}\Sigma^{1/2\top}) = \mathcal{N}(\mu, \Sigma).$$

4.2 Inference

To perform inference using the variational approximation, we need to compute the integral over all models θ ,¹⁸

$$\begin{aligned} p(y^* | \mathbf{x}^*, \mathbf{y}) &= \int p(y^* | \mathbf{x}^*, \theta) p(\theta | \mathbf{y}) d\theta \\ &\approx \int p(y^* | \mathbf{x}^*, \theta) q_\lambda(\theta) d\theta \\ &= \int \int p(y^* | \mathbf{x}^*, \theta, f^*) q_\lambda(\theta | f^*) p(f^* | \mathbf{x}^*, \theta) d\theta df^* \\ &= \int \int p(y^* | f^*) p(f^* | \mathbf{x}^*, \theta) q_\lambda(\theta) d\theta df^* \\ &= \int p(y^* | f^*) \int p(f^* | \mathbf{x}^*, \theta) q_\lambda(\theta) d\theta df^* \\ &= \int p(y^* | f^*) q_\lambda(f^* | \mathbf{x}^*) df^*. \end{aligned}$$

¹⁸ The key insight here is that parameters θ do not matter, only their results f^* .

Marginalize over f^*

Marginalize out θ

Thus, we have reduced the high-dimensional integral over the parameters θ to a one-dimensional integral over f^* . While this integral is generally still intractable, it can be approximated efficiently.

5 Markov chain Monte Carlo

Markov chain Monte Carlo (MCMC) methods seek to approximate the intractable distribution p by drawing m approximate samples from p ,

$$\begin{aligned} p(y^* | x^*, y) &= \int p(y^* | x^*, \theta) p(\theta | y) d\theta \\ &= \mathbb{E}_{\theta \sim p(\cdot | y)} [p(y^* | x^*, \theta)] \\ &\approx \frac{1}{m} \sum_{i=1}^m p(y^* | x^*, \theta^{(i)}). \end{aligned} \quad \text{Law of large numbers}$$

Using *Hoeffding's inequality*, we can compute a bound on the error,

$$p \left(\left| \mathbb{E}_{\theta \sim p(\cdot | y)} [p(y^* | x^*, \theta)] - \frac{1}{m} \sum_{i=1}^m p(y^* | x^*, \theta^{(i)}) \right| > \epsilon \right) \leq 2 \exp(-2m\epsilon^2).$$

Thus, the probability of error decreases exponentially in m . To get a probability $\leq \delta$ of error $> \epsilon$, we need

$$\begin{aligned} 2 \exp(-2m\epsilon^2) &\leq \delta \\ -2m\epsilon^2 &\leq \log \frac{\delta}{2} \\ m &\geq \frac{\log 2 - \log \delta}{2\epsilon^2} \end{aligned}$$

samples. Intuitively, if we want a lower error or lower probability, we need more samples.

However, we cannot sample directly from the posterior $p(\theta | y)$, because it is intractable. The key idea of MCMC is to construct a Markov chain with stationary distribution $p(\theta | y)$ and sample from that.

5.1 Markov chains

Definition 5.1 (Markov chain). A Markov chain is a sequence of random variables $(X_t)_{t \in \mathbb{N}_0}$ with prior $P(X_1)$ and transition probabilities $P(X_{t+1} | X_t)$ independent of t . The Markov assumption is thus the following,

$$X_{t+1} \perp X_{1:t-1} | X_t.$$

Intuitively, this states that future behavior is independent of past states given the present state. In other words, all past information is encapsulated by the current state.

Definition 5.2 (Stationary distribution). A distribution π is stationary with respect to the transition function P iff $P(X_n) = P(X_{n+1})$. In other words, the probability distribution over states remains the same between timesteps. The following must hold for all x ,

$$\pi(x) = \sum_{x'} P(x | x') \pi(x').$$

Definition 5.3 (Ergodicity). A Markov chain is ergodic iff there exists a $t \in \mathbb{N}_0$ such that for any x, x' , the following holds,

$$P^{(t)}(x' | x) > 0,$$

where $P^{(t)}(x' | x)$ is the probability to reach x' from x in *exactly* t steps. Intuitively, this means that any state is reachable from another within the same amount of steps.

Remark. An easy way of ensuring that a Markov chain is ergodic is to add “self-loops” to every vertex.

Theorem 5.4 (Fundamental theorem of ergodic Markov chains). An ergodic Markov chain has a unique and positive stationary distribution $\pi(X) > 0$ such that for all x , the following holds,

$$\lim_{n \rightarrow \infty} P(X_n = x) = \pi(x),$$

independent of the initial distribution $P(X_1)$. *I.e.*, ergodic Markov chains always converge to a stationary distribution.

Thus, making use of the fundamental theorem of ergodic Markov chains, we can construct an ergodic Markov chain such that its stationary distribution coincides with the posterior distribution, $\pi(x) = p(\theta | \mathbf{y})$. If we then sample “sufficiently long” from this Markov chain, X_t is drawn from a distribution that is “very close” to the stationary distribution π , which is equal to the posterior distribution.

Definition 5.5 (Detailed balance equation). A Markov chain satisfies the detailed balance equation for an unnormalized distribution q iff the following holds for any x, x' ,

$$q(x)P(x' | x) = q(x')P(x | x').$$

Theorem 5.6. If a finite Markov chain satisfies the detailed balance equation with respect to q , then $\frac{1}{Z}q$ is a stationary distribution.

Proof. Let $p_t = q$. Then for any x , the following holds,

$$\begin{aligned}
 p_{t+1}(x) &= \sum_{x'} P(x | x') p_t(x') && \text{Markov assumption} \\
 &= \sum_{x'} P(x | x') q(x') \\
 &= \sum_{x'} P(x' | x) q(x) && \text{Detailed balance equation} \\
 &= q(x) \sum_{x'} P(x' | x) \\
 &= q(x).
 \end{aligned}$$

Thus, q is the stationary distribution. ■

5.2 Sampling

If we can show that the detailed balance equation holds for the *unnormalized* posterior distribution, then we know that the posterior distribution is the stationary distribution of the Markov chain.¹⁹ Thus, we do not need to know the true posterior. It suffices to know its unnormalized version.²⁰

The *Metropolis-Hastings algorithm* constructs a Markov chain with the posterior as stationary distribution. It uses an arbitrary proposal transition distribution $R(x' | x)$ which, given we are in state x , proposes a new state x' .²¹ Following the proposal with probability

$$\alpha(x' | x) \doteq \min \left\{ 1, \frac{q(x') R(x | x')}{q(x) R(x' | x)} \right\},$$

yields a Markov chain with the desired stationary distribution $\frac{1}{Z} q(x)$, because it makes it satisfy the detailed balance equation.

¹⁹ The only problem is that we do not know the rate of convergence to the stationary distribution.

²⁰ Recall that the normalizer was the intractable part of $p(\theta | y)$.

²¹ The rate of convergence of this algorithm strongly depends on the choice of R .

```

1: function METROPOLISHASTINGS(R)
2:   initialize  $x$ 
3:   for  $t = 1, \dots, T$  do
4:      $x' \sim R(x' | x)$ 
5:      $u \sim \text{Unif}([0, 1])$ 
6:     if  $u \leq \alpha(x' | x)$  then
7:        $x \leftarrow x'$ 
8:     end if
9:   end for
10: end function
    
```

Algorithm 1. The Metropolis-Hastings algorithm. Each iteration, with a random probability, follow the proposal distribution.

5.3 Proposal distributions

We want to converge to the stationary distribution as fast as possible. The proposal distribution has a big influence on this.

Gibbs sampling. A popular example algorithm for specifying a proposal distribution R is Gibbs sampling. Gibbs sampling works by iteratively improving the variables. It starts with an initial assignment x to all variables, fixing the observed variables to their observed value. Then, iteratively uniformly pick a variable X_i to update given the rest of the set values by sampling $p(X_i \mid x_{1:i-1}, x_{i+1:n})$.²²

²² Sampling from this distribution is typically efficient.

```

1: function GIBBSAMPLING
2:   initialize  $\mathbf{x} = [x_1, \dots, x_n] \in \mathbb{R}^n$ 
3:   for  $t = 1, \dots, T$  do
4:     uniformly sample  $i$  from  $\{1, \dots, n\}$ 
5:      $\mathbf{x}_{-i} \leftarrow [x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n]$ 
6:     update  $x_i$  by sampling from  $p(x_i \mid \mathbf{x}_{-i})$ 
7:   end for
8: end function

```

Algorithm 2. Gibbs sampling.

Then, Gibbs sampling is a Metropolis-Hastings algorithm with the following proposal distribution,

$$R(x' \mid x) \doteq \begin{cases} p(x'_i \mid x'_{1:i-1}, x'_{i+1:n}) & x' \text{ differs from } x \\ 0 & \text{otherwise} \end{cases}$$

and acceptance distribution $\alpha(x' \mid x) = 1$ for all x, x' .

Gaussian. Generally, we focus on positive distributions written as the following,

$$p(x) = \frac{1}{Z} \exp(-f(x)),$$

where f is called an energy function (high energy \equiv low probability, low energy \equiv high probability). Then, the acceptance distribution becomes the following,

$$\alpha(x' \mid x) = \min \left\{ 1, \frac{R(x \mid x')}{R(x' \mid x)} \exp(f(x) - f(x')) \right\}.$$

One option for the proposal distribution is

$$R(x' \mid x) = \mathcal{N}(x'; x, \tau \mathbf{I}).$$

Since this R is symmetric,

$$\frac{R(x \mid x')}{R(x' \mid x)} = 1.$$

Thus, if R proposes to move to a region with lower energy, the acceptance probability will always be 1. If R proposes to move to a region with higher energy, the probability moves toward 0 dependent on how much higher the region is.

However, we want to move as quickly as possible through the function, going through all high-density areas. But, this R is “uninformed” and

thus proposes to go in any direction. We would like to propose areas with lower energy, which are high-density areas for $p(x)$. Then, we will have less iterations where the proposal simply gets rejected.

Metropolis adjusted Langevin algorithm. An improvement to the Gaussian proposal distribution is the Metropolis adjusted Langevin algorithm (MALA),

$$R(x' | x) = \mathcal{N}(x'; x - \tau \nabla f(x), 2\tau \mathbf{I}).$$

This proposes moving to high-density areas, thus it much more efficiently converges to the stationary distribution.

The problem with this is that it requires access to the full energy function to compute the gradient, which can be expensive for large datasets. This is solved by stochastically estimating the gradient $\nabla f(x)$.

6 Bayesian neural networks

So far, we have explored techniques for computing uncertainty of linear models.²³ However, in practice, we can often get better performance by considering non-linear dependencies. Thus, this is what we will explore next.

Neural networks typically look like the following,

$$f_{\theta}(\mathbf{x}) = \phi(\mathbf{W}_{\ell}\phi(\mathbf{W}_{\ell-1}\cdots\phi(\mathbf{W}_1\mathbf{x}))).$$

Bayesian neural network models specify a prior distribution over the weights,

$$\theta \sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I}),$$

and use likelihood distributions parametrized by a neural network,

$$y \mid \mathbf{x}, \theta \sim \mathcal{N}(f(\mathbf{x}; \theta), \sigma^2),$$

which assumes homoscedastic noise.²⁴ However, we can also parameterize a likelihood that can model heteroscedastic noise by predicting the variance of the Gaussian,

$$y \mid \mathbf{x}, \theta \sim \mathcal{N}(f_{\mu}(\mathbf{x}; \theta), \exp f_{\sigma^2}(\mathbf{x}; \theta)).$$

The MAP estimate of a BNN is the following,

$$\begin{aligned} \hat{\theta} &= \underset{\theta}{\operatorname{argmax}} p(\theta \mid \mathbf{X}, \mathbf{y}) \\ &= \underset{\theta}{\operatorname{argmax}} p(\theta) p(\mathbf{y} \mid \mathbf{X}, \theta) \\ &= \underset{\theta}{\operatorname{argmin}} -\log p(\theta) - \sum_{i=1}^n \log p(y_i \mid \mathbf{x}_i, \theta) \\ &= \underset{\theta}{\operatorname{argmin}} -\lambda \|\theta\|^2 + \sum_{i=1}^n -\log \mathcal{N}(y_i; f_{\mu}(\mathbf{x}_i; \theta), f_{\sigma^2}(\mathbf{x}_i; \theta)) \\ &= \underset{\theta}{\operatorname{argmin}} -\lambda \|\theta\|^2 + \sum_{i=1}^n -\log \left(\frac{1}{2\pi f_{\sigma^2}(\mathbf{x}_i; \theta)} \exp \left(-\frac{1}{2f_{\sigma^2}(\mathbf{x}_i; \theta)} (y_i - f_{\mu}(\mathbf{x}_i; \theta))^2 \right) \right) \\ &= \underset{\theta}{\operatorname{argmin}} -\lambda \|\theta\|^2 + \sum_{i=1}^n \log(2\pi) + \log(f_{\sigma^2}(\mathbf{x}_i; \theta)) + \frac{1}{2f_{\sigma^2}(\mathbf{x}_i; \theta)} (y_i - f_{\mu}(\mathbf{x}_i; \theta))^2 \\ &= \underset{\theta}{\operatorname{argmin}} -\lambda \|\theta\|^2 + \sum_{i=1}^n \log(f_{\sigma^2}(\mathbf{x}_i; \theta)) + \frac{(y_i - f_{\mu}(\mathbf{x}_i; \theta))^2}{2f_{\sigma^2}(\mathbf{x}_i; \theta)}. \end{aligned}$$

Thus, the MAP estimate is a balance between the mean and variance predictions. If we perfectly predict y_i with f_{μ} , we only need to make f_{σ^2} smaller. Otherwise, we can attenuate for the error $(y_i - f_{\mu}(\mathbf{x}_i; \theta))^2$ with f_{σ^2} in its denominator, for which we have to pay logarithmically. Intuitively, the model can attenuate certain losses for certain datapoints by attributing the error to large variance.

However, the problem with the MAP estimate is that it does not use the entire distribution over θ . In other words, it does not account for epistemic uncertainty, only aleatoric. But, just like before, using the whole

²³ The likelihood have parameters linearly dependent on the input feature.

²⁴ Same noise for all data points.

distribution would be intractable. Thus, we need some approximation techniques to make it tractable. We will explore this in the next subsections.

6.1 Variational inference

Since BNNs are just a distribution over the weights θ , we can approximate its distribution with variational inference. Then, we can learn the distribution over θ by optimizing the ELBO of the approximation distribution q_λ . Then, we can do inference as follows,

$$\begin{aligned}
 p(y^* | x^*, X, y) &= \int p(y^* | x^*, \theta) p(\theta | X, y) d\theta \\
 &= \mathbb{E}_{\theta \sim p(\cdot | X, y)} [p(y^* | x^*, \theta)] \\
 &\approx \mathbb{E}_{\theta \sim q_\lambda} [p(y^*, x^*, \theta)] && \text{Variational inference} \\
 &\approx \frac{1}{m} \sum_{j=1}^m p(y^* | x^*, \theta^{(j)}), \quad \theta^{(j)} \sim q_\lambda. && \text{Monte Carlo}
 \end{aligned}$$

If q_λ is Gaussian, then the approximate predictive distribution becomes a mixture of Gaussians. The mean and variance of this distribution are the following,

$$\begin{aligned}
 \mathbb{E}[y^* | x^*, X, y] &\approx \bar{\mu}(x^*) \doteq \frac{1}{m} \sum_{j=1}^m f_\mu(x^*; \theta^{(j)}) \\
 \text{Var}[y^* | x^*, X, y] &= \mathbb{E}_\theta [\text{Var}_{y^*}[y^* | x^*, \theta]] + \text{Var}_\theta [\mathbb{E}_{y^*}[y^* | x^*, \theta]] && \text{Law of total variance} \\
 &\approx \underbrace{\frac{1}{m} \sum_{j=1}^m f_{\sigma^2}(x^*; \theta^{(j)})}_{\text{aleatoric uncertainty}} + \underbrace{\frac{1}{m-1} \sum_{j=1}^m \left(f_\mu(x^*; \theta^{(j)}) - \bar{\mu}(x^*) \right)^2}_{\text{epistemic uncertainty}}
 \end{aligned}$$

6.2 Markov chain Monte Carlo

It is also possible to apply MCMC to BNNs. MCMC methods produce a sequence of weights $\theta^{(1)}, \dots, \theta^{(T)}$. Using the ergodic theorem we can then make predictions with the following,

$$p(y^* | x^*, X, y) \approx \frac{1}{T} \sum_{j=1}^T p(y^* | x^*, \theta^{(j)}).$$

However, models are often very large, so we cannot store T times the parameters of the network $\mathcal{O}(Td)$. Thus, we need to approximate. A simple solution is to only keep a subset of m weights. But, we can also approximate the distribution with a Gaussian,

$$\theta \sim \mathcal{N}(\mu, \Sigma),$$

and keeping running averages, only requiring $\mathcal{O}(d^2)$ space complexity, where

$$\mu = \frac{1}{T} \sum_{j=1}^T \theta^{(j)}, \quad \Sigma = \frac{1}{T-1} \sum_{j=1}^T (\theta^{(j)} - \mu)(\theta^{(j)} - \mu)^\top.$$

SWAG [Maddox et al., 2019] is an example of a model that does this, but instead of an MCMC method, it uses stochastic gradient descent to sample models.

6.3 Monte Carlo dropout

Dropout regularization is often used in traditional neural networks to improve generalization. It works by randomly selecting weights to set to 0. However, using *Monte Carlo dropout*, we can view this as performing variational inference. Let p be the probability that we omit parameter, then the variational posterior is given as the following,

$$q(\boldsymbol{\theta} \mid \lambda) = \prod_{j=1}^d q_j(\theta_j \mid \lambda_j)$$

$$q_j(\theta_j \mid \lambda_j) = p\delta_0(\theta_j) + (1 - p)\delta_{\lambda_j}(\theta_j),$$

where d is the number of parameters in the neural network. Intuitively, this posterior says that the j -th weight has value 0 with probability p and value λ_j with probability $1 - p$.

The difference with dropout regularization is that we also need to use dropout during inference for this to be variational inference,

$$p(y^* \mid \mathbf{x}^*, \mathbf{y}) \approx \mathbb{E}_{\boldsymbol{\theta} \sim q_\lambda} [p(y^* \mid \mathbf{x}^*, \boldsymbol{\theta})]$$

$$\approx \frac{1}{m} \sum_{j=1}^m p(y^* \mid \mathbf{x}^*, \boldsymbol{\theta}^{(j)})$$

$$\boldsymbol{\theta}^{(j)} \stackrel{i.i.d.}{\sim} q_\lambda.$$

Intuitively, we average the distribution of m neural networks for each of which we randomly drop out weights.

6.4 Probabilistic ensembles

We have seen that variational inference can be seen as averaging the predictions of m neural networks. A natural adaptation of this idea is to learn the weights of m neural networks. The idea is to randomly choose m training subsets, each with n data points. Then, we compute m MAP estimates $\boldsymbol{\theta}^{(j)}$, yielding the following approximation,

$$p(y^* \mid \mathbf{x}^*, \mathbf{y}) \approx \frac{1}{m} \sum_{j=1}^m p(y^* \mid \mathbf{x}^*, \boldsymbol{\theta}^{(j)}).$$

6.5 Calibration

A key challenge of BNNs is *calibration*. We want models to be well-calibrated, which means that the confidence that they have in their predictions coincides with the accuracy they have over many samples. For example, let's say that we have a classification model that predicts that a

data point belongs to a certain class with 80% probability. If the model is well-calibrated, then the prediction should be correct 80% of the time. We can calibrate models by adjusting the probability estimation of models.

Reliability diagrams (Figure 6.1) are a way of determining how calibrated a model is. This diagram is constructed by making predictions on a validation dataset. These predictions are then divided into M bins according to the predicted class probability,

$$B_j = \left\{ y \mid p(y) \in \left[\frac{j-1}{m}, \frac{j}{m} \right) \right\}.$$

Within each bin, we then compare the predicted probabilities (confidence) with how often the input actually belonged to the class (frequency),

$$\begin{aligned} \text{conf}(B_j) &= \frac{1}{|B_j|} \sum_{y \in B_j} p_{\theta}(y) \\ \text{acc}(B_j) &= \frac{1}{|B_j|} \sum_{y \in B_j} \mathbb{1}\{y = \hat{y}\}, \end{aligned}$$

A model is well-calibrated if $\text{conf}(B_j) \approx \text{acc}(B_j)$ for all bins B_j .



Figure 6.1. Reliability diagrams. The top diagram shows a well-calibrated model, the second diagram is overconfident, and the third diagram is underconfident.

7 Active learning

Until now, we have only covered how to learn/represent aleatoric and epistemic uncertainty in machine learning. *Active learning* covers how to use this measure for deciding which data to collect. Intuitively, we want to collect data points in places where we would gain the most information, *i.e.*, where the uncertainty is high. This assumes that there is some cost associated with collecting data points, making it important to be careful about which data points to pick that maximize the information obtained.

Let \mathcal{X} be a set of possible observations of f , and y_x the observation at $x \in \mathcal{X}$,

$$y_x \doteq f(x) + \epsilon_x, \quad \epsilon_x \in \mathcal{N}(\mathbf{0}, \sigma_n^2 \mathbf{I}).$$

Then, we want to observe a subset $S \subseteq \mathcal{X}$ of a fixed size that maximizes the information gain between the model f and y_S , which yields the following maximization objective,

$$I(S) \doteq I(f_S; y_S) = H[f_S] - H(f_S | y_S),$$

where $H[f_S]$ denotes the uncertainty about f_S before obtaining the observations y_S and $H[f_S | y_S]$ corresponds to the uncertainty about f_S after obtaining observations y_S . This problem is \mathcal{NP} -hard, thus we need to formulate a strategy.

7.1 Sampling strategies

Uncertainty sampling. The simplest strategy is to greedily pick points one by one, which entails that we pick the locations, x_1, \dots, x_n , individually by greedily finding the location with maximal mutual information. That is, if we have already picked locations $S_t = \{x_1, \dots, x_t\}$, then the next point, x_{t+1} , maximizes its mutual information,

$$x_{t+1} \doteq \operatorname{argmax}_{x \in \mathcal{X}} I(f_x; y_x | y_{S_t}).$$

Assuming that f is modeled by a Gaussian,

$$\begin{aligned} x_{t+1} &= \operatorname{argmax}_{x \in \mathcal{X}} \frac{1}{2} \log \left(1 + \frac{\sigma_{x|S_t}^2}{\sigma_n^2(x)} \right) \\ &= \operatorname{argmax}_{x \in \mathcal{X}} \frac{\sigma_{x|S_t}^2}{\sigma_n^2(x)}. \end{aligned}$$

Assuming that the label noise is independent of x , *i.e.*, homoscedastic,

$$x_{t+1} = \operatorname{argmax}_{x \in \mathcal{X}} \sigma_{x|S_t}^2.$$

Thus, if f is modeled by a Gaussian and we assume homoscedastic noise, greedily maximizing mutual information corresponds to picking the point x with the largest variance.

Due to the *information never hurts* principle, mutual information is monotone submodular, which means that adding data points can only increase the mutual information in expectation. From this, it follows that the greedy algorithm provides a constant-factor approximation, which means that uncertainty sampling is near-optimal.

Heteroscedastic noise. If the data is not homoscedastic, uncertainty sampling will fail, because it fails to distinguish between epistemic and aleatoric uncertainty. Therefore, the most uncertain point is not necessarily the most informative one. Thus, maximizing the mutual information yields the following,

$$\mathbf{x}_{t+1} = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \frac{\sigma_{x|S_t}^2}{\sigma_h^2(\mathbf{x})}.$$

Here we make a trade-off between large epistemic uncertainty and large aleatoric uncertainty. Ideally, we find an x where the epistemic uncertainty is large, and the aleatoric uncertainty low, because we want to make sure we get a lot of information (high epistemic uncertainty), but also that we can have high confidence in the point we choose (low aleatoric uncertainty).

Bayesian active learning by disagreement. Uncertainty sampling in classification corresponds to selecting samples that maximize entropy of the predicted label, *i.e.*, points that are close to the decision boundary. However, the uncertainty in points around the decision boundary are often due to aleatoric noise. Hence, we will not learn much from observing points there.

Thus, we need to distinguish between aleatoric and epistemic uncertainty of f_θ ,

$$\begin{aligned} \mathbf{x}_{t+1} &= \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} I(y_x; \theta \mid \mathbf{x}_{1:t}, y_{1:t}) \\ &= \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} H[y_x \mid \mathbf{x}_{1:t}, y_{1:t}] - H[y_x \mid \theta, \mathbf{x}_{1:t}, y_{1:t}] \\ &= \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \underbrace{H[y_x \mid \mathbf{x}_{1:t}, y_{1:t}]}_{\text{entropy of pred. posterior}} - \underbrace{\mathbb{E}_{\theta \mid \mathbf{x}_{1:t}, y_{1:t}}[H[y_x \mid \theta]]}_{\text{entropy of likelihood}}. \end{aligned}$$

The first term measures the entropy of the average prediction, while the second term measures the average entropy of predictions. Thus, the first term looks for points where the average prediction is uncertain. In contrast, the second term penalizes points where many of the sampled models θ are uncertain about their prediction. Thus, we want to find points x where the posterior is uncertain (epistemic uncertainty) because of all models θ being extremely certain about their differing predictions (aleatoric uncertainty).

Since $p(\theta \mid \mathbf{x}_{1:t}, y_{1:t})$ is intractable, we need to use variational inference

and Monte Carlo to approximate the second term,

$$\begin{aligned}\mathbb{E}_{\boldsymbol{\theta}|x_{1:t}, y_{1:t}}[H[y_x | \boldsymbol{\theta}]] &\approx \mathbb{E}_{q_\lambda}[H[y_x | \boldsymbol{\theta}]] && \text{Variational inference} \\ &\approx \frac{1}{m} \sum_{i=1}^m H[y_x | \boldsymbol{\theta}^{(i)}], \quad \boldsymbol{\theta}^{(i)} \sim q_\lambda. && \text{Monte Carlo}\end{aligned}$$

We can use another approximation method, such as variational inference, Markov chain Monte Carlo, or SWAG, to approximate the predictive posterior in the first term.

8 Bayesian optimization

In *Bayesian optimization*, we do not only want to reduce uncertainty, but we also want to maximize some objective.²⁵ This means that we have to make a trade-off between exploration (minimizing uncertainty) and exploitation (maximizing performance).

Definition 8.1 (Regret). The cumulative regret for a time horizon T associated with choices x_1, \dots, x_T is defined as

$$R_T \doteq \sum_{t=1}^T \underbrace{\left(\max_x f(x) - f(x_t) \right)}_{\text{"instantaneous regret"}}.$$

The regret can be interpreted as the additive loss with respect to the maximally achievable value $\max_x f(x)$.

The goal is to find algorithms that achieve a sublinear regret,

$$\lim_{T \rightarrow \infty} \frac{R_T}{T} = 0.$$

This leads to the algorithm converging on the maximum $\max_x f(x)$. Note that using an algorithm that explores forever will result in the regret growing linearly, because we will never settle on a maximum value. In contrast, if we use an algorithm that never explores and thus only exploits, we might never find $\max_x f(x)$. Thus, achieving sublinear regret requires balancing exploration and exploitation.

8.1 Acquisition functions

Upper confidence bound. The principle of *optimism in the face of uncertainty* naturally suggests picking the point where we can hope for the optimal outcome. This corresponds to maximizing the *upper confidence bound* (UCB),²⁶

$$x_t = \operatorname{argmax}_{x \in \mathcal{X}} \mu_{t-1}(x) + \beta_t \sigma_{t-1}(x),$$

where β_t regulates how confident we are in our current model. A high β_t leads the model to explore, while a low β_t leads the model to exploit. This acquisition function naturally trades exploitation by preferring a large posterior mean with exploration by preferring a large posterior variance.

If f can be represented by our model and we choose β_t “correctly”,

$$R_T \in \mathcal{O}^*\left(\sqrt{\frac{\gamma_T}{T}}\right),$$

where

$$\gamma_T = \max_{|S| \leq T} I(f; y_S),$$



Figure 8.1. Illustration of Bayesian optimization. We pass an input x_t into the unknown function f^* to obtain a noisy observation y_t .

²⁵ This is in contrast with active learning. An example of this is automatically tuning the hyperparameters of a machine learning model. In this scenario, we would both want to minimize the uncertainty of the parameter space and maximize the performance of the eventual model.

²⁶ This assumes that the parameters of the model (e.g. GP) are known. However, this is not the case in practice. Thus, in practice, we can alternate between learning the hyperparameters from the currently selected data, and selecting the next data-point. However, this introduces a danger of overfitting. This can be solved by either placing a hyperprior on the hyperparameters, or occasionally selecting some points at random.

quantifies the maximum information gain. Thus, the maximum information gain γ_T determines the regret of the UCB algorithm.

Thompson sampling. At every iteration t , Thompson sampling draws a function realization from the Gaussian process,

$$\tilde{f}_t \sim p(f \mid \mathbf{x}_{1:t}, \mathbf{y}_{1:t}),$$

and selects

$$\mathbf{x}_{t+1} = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \tilde{f}_t(\mathbf{x}).$$

The randomness in the realization of \tilde{f}_t is sufficient to trade exploration and exploitation. Similarly to UCB, it has sublinear regret bounds.

Kernel	γ_T bound
Linear	$\mathcal{O}(d \log T)$
Gaussian	$\mathcal{O}((\log T)^{d+1})$
Matérn ($\nu > 1/2$)	$\mathcal{O}\left(T^{\frac{d}{2\nu+d}} (\log T)^{\frac{2\nu}{2\nu+d}}\right)$

Table 1. Information gain bounds of common Gaussian process kernels. These guarantee sublinear regret, which means that they are guaranteed to converge to the maximum value of the function.

9 Markov decision processes

Definition 9.1 (Markov decision process). A Markov decision process $\langle \mathcal{X}, \mathcal{A}, p, r \rangle$ is specified by the following,

- Set of states \mathcal{X} (not necessarily finite);
- Set of actions \mathcal{A} (not necessarily finite);
- Initial state distribution $p(x_0)$;
- Transition probabilities $p(x' | x, a)$;
- Reward function $r(x, a)$ (or $r(x, a, x')$).

As can be seen by the form of the transition probabilities, we make the Markov assumption, where the future state only depends on the current state and action.

In general, we want to maximize the long-term reward, either over a finite horizon T ,

$$\max \mathbb{E} \left[\sum_{t=0}^T r(x_t, a_t) \right],$$

or over an infinite horizon where future rewards are discounted by a decaying factor $\gamma \in (0, 1)$,²⁷

$$\max \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(x_t, a_t) \right].$$

A *policy* $\pi : \mathcal{X} \rightarrow \mathcal{A}$ is a function that maps states to their action. A policy π induces a Markov chain with $P(X_{t+1} = x' | X_t = x) = p(x' | x, \pi(x))$ as transition probabilities. In general, we want to find the policy that maximizes the expected value,

$$J(\pi) \doteq \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(X_t, \pi(X_t)) \right].$$

9.1 Bellman expectation equation

The *value function* of a state x is the expected sum of discounted future rewards, obtained from subsequent states. This is defined as the following,

$$V^\pi(x) \doteq \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(X_t, \pi(X_t)) \mid X_0 = x \right],$$

For now, we will assume that p and r are known, but, as we will see in section 10, in reinforcement learning, these are unknown.

²⁷ We cannot maximize the summed future reward (without decaying factor), because this could be ∞ for all possible strategies.

We can solve for the value function of a policy using a bit of linear algebra. Let $\mathcal{X} = \{1, \dots, n\}$, then we can define the following vectors and matrix,

$$\begin{aligned} \mathbf{v}^\pi &= \begin{bmatrix} V^\pi(1) \\ \vdots \\ V^\pi(n) \end{bmatrix} \\ \mathbf{r}^\pi &= \begin{bmatrix} r(1, \pi(1)) \\ \vdots \\ r(n, \pi(n)) \end{bmatrix} \\ \mathbf{T}^\pi &= \begin{bmatrix} p(1 | 1, \pi(1)) & \cdots & p(n | 1, \pi(1)) \\ \vdots & \ddots & \vdots \\ p(1 | n, \pi(n)) & \cdots & p(n | n, \pi(n)) \end{bmatrix}. \end{aligned}$$

Using the Bellman expectation equation, we have the following equality,

$$\begin{aligned} \mathbf{v}^\pi &= \mathbf{r}^\pi + \gamma \mathbf{T}^\pi \mathbf{v}^\pi \\ \iff (\mathbf{I} - \gamma \mathbf{T}^\pi) \mathbf{v}^\pi &= \mathbf{r}^\pi \\ \iff \mathbf{v}^\pi &= (\mathbf{I} - \gamma \mathbf{T}^\pi)^{-1} \mathbf{r}^\pi. \end{aligned}$$

However, this is computationally expensive. We could also use fixed-point iteration to obtain an (approximate) solution. This involves iteratively computing the value function of each node, given the last value function,

$$\mathbf{v}_t^\pi = \mathbf{r}^\pi + \gamma \mathbf{T}^\pi \mathbf{v}_{t-1}^\pi.$$

This is faster, because \mathbf{T}^π is sparse.

which contains a recursive relationship,

$$\begin{aligned}
V^\pi(x) &\doteq \mathbb{E}_x \left[\sum_{t=0}^{\infty} \gamma^t r(X_t, \pi(X_t)) \mid X_0 = x \right] \\
&= \mathbb{E}_x \left[r(X_0, \pi(X_0)) + \sum_{t=1}^{\infty} \gamma^t r(X_t, \pi(X_t)) \mid X_0 = x \right] \\
&= \mathbb{E}_x[r(X_0, \pi(X_0)) \mid X_0 = x] + \mathbb{E}_x \left[\sum_{t=0}^{\infty} \gamma^{t+1} r(X_{t+1}, \pi(X_{t+1})) \mid X_0 = x \right] \\
&= r(x, \pi(x)) + \gamma \mathbb{E}_{x'} \left[\mathbb{E}_x \left[\sum_{t=0}^{\infty} \gamma^t r(X_{t+1}, \pi(X_{t+1})) \mid X_1 = x' \right] \right] && \text{Condition on realization of } X_1 \\
&= r(x, \pi(x)) + \gamma \sum_{x' \in \mathcal{X}} p(x' \mid x, \pi(x)) \cdot \mathbb{E}_{x'} \left[\sum_{t=0}^{\infty} \gamma^t r(X_{t+1}, \pi(X_{t+1})) \mid X_1 = x' \right] \\
&= r(x, \pi(x)) + \gamma \sum_{x' \in \mathcal{X}} p(x' \mid x, \pi(x)) \cdot \mathbb{E}_{x'} \left[\sum_{t=0}^{\infty} \gamma^t r(X_t, \pi(X_t)) \mid X_0 = x' \right] && \text{Stationarity of Markov chains} \\
&= r(x, \pi(x)) + \gamma \sum_{x' \in \mathcal{X}} p(x' \mid x, \pi(x)) \cdot V^\pi(x').
\end{aligned}$$

This equation is known as the *Bellman expectation equation*. Intuitively, this means that the value of the current state corresponds to the reward from the next action, plus the discounted sum of expected future rewards obtained from the subsequent states (which are their values).

Theorem 9.2 (Bellman's theorem). A policy π^* is optimal iff it is greedy w.r.t. its own value function V^* .

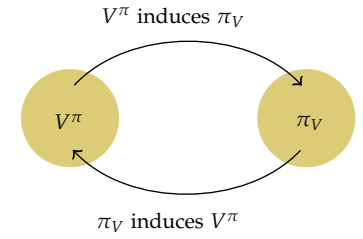


Figure 9.1. Cyclic dependency between value function and greedy policy.

9.2 Policy iteration

Every value function induces a policy where we greedily choose the action that maximizes the expected value,

$$\pi_V(x) = \operatorname{argmax}_{a \in \mathcal{A}} r(x, a) + \gamma \sum_{x' \in \mathcal{X}} p(x' \mid x, a) V(x'), \quad (1)$$

while the policy induces a value function as we have seen. This leads us to *policy iteration*, where we find an optimal policy by alternating between computing the value function w.r.t. π (using fixed-point iteration) and computing the next greedy policy w.r.t. V^π (using Equation (1)), until convergence.

Policy iteration is expensive, because every iteration requires computing a value function, but it is guaranteed to converge monotonically.

9.3 Value iteration

Value iteration can be seen as a dynamic program that computes the optimal value function, where the state is how many timesteps t we look

```

function POLICYITERATION( $\langle \mathcal{X}, \mathcal{A}, p, r \rangle$ )
  randomly initialize  $\pi$ 
  while not converged do
     $V \leftarrow \text{VALUEFUNCTION}(\pi)$ 
     $\pi \leftarrow \text{GREEDYPOLICY}(V)$ 
  end while
  return  $\pi$ 
end function

```

Algorithm 3. Policy iteration algorithm that finds an exact solution in a polynomial number of iterations.

ahead, and $V_t(x)$ is the value function of that. The recurrence relationship is then the following,

$$V_0(x) = \max_{a \in \mathcal{A}} r(x, a)$$

$$V_t(x) = \max_{a \in \mathcal{A}} r(x, a) + \gamma \sum_{x' \in \mathcal{X}} p(x' | x, a) V_{t-1}(x').$$

Since V_{t-1} looks $t-1$ steps into the future, V_t looks t steps ahead. We keep iterating until ϵ -optimal convergence.²⁸

²⁸ ϵ -optimal in the sense that the largest difference between $V^*(x)$ and $V_t(x)$ for any x is at most ϵ .

```

function VALUEITERATION( $\langle \mathcal{X}, \mathcal{A}, p, r \rangle$ )
  for  $x \in \mathcal{X}$  do                                      $\triangleright$  Initialize  $V_0$ 
     $V_0(x) \leftarrow \max_{a \in \mathcal{A}} r(x, a)$ 
  end for
   $t \leftarrow 0$ 
  while  $\|v_t - v_{t-1}\|_\infty > \epsilon$  do
     $t \leftarrow t + 1$                                  $\triangleright$  Look one more step into the future
    for  $x \in \mathcal{X}$  do
       $V_t(x) \leftarrow \max_{a \in \mathcal{A}} r(x, a) + \gamma \sum_{x' \in \mathcal{X}} p(x' | x, a) V_{t-1}(x')$ 
    end for
  end while
  return  $\text{GREEDYPOLICY}(V_t)$ 
end function

```

Algorithm 4. Value iteration algorithm that finds an ϵ -optimal solution in a polynomial number of iterations. v_t is the vector containing all values of V_t . $\|x\|_\infty$ is the largest value of x .

Recall Theorem 9.2 that states that if a policy is greedy w.r.t. its own value function, it is optimal. Thus, after finding an ϵ -optimal value function, we can simply choose the policy that greedily picks according to this value function.

Value iteration is not guaranteed to converge monotonically, but it is guaranteed to converge to an ϵ -optimal policy in polynomial time. Furthermore, value iteration is inexpensive, compared to policy iteration.

10 Reinforcement learning

In *reinforcement learning* (RL), we are concerned with acting in unknown environments. These environments are still modeled by MDPs, but in RL, we do not have access to the transition probabilities p and reward function r . Thus, RL is at the intersection of probabilistic planning (MDPs) and learning, *i.e.*, everything we have learned thus far comes together here.

Remark. We will start by assuming that the state-action space is finite. Then, we will move on to potentially infinite state spaces. After that, we will learn about infinite action spaces.

Since the environment is unknown, we need to explore the state-action space to find where the reward lies. However, we also want to act optimally by exploiting what we have learned thus far. This is called the exploration/exploitation dilemma and is what algorithms need to solve.

Another way that reinforcement learning differs from supervised learning is that data depends on past actions. Trajectory data looks like the following,

$$\tau = (\langle x_0, a_0, r_0, x_1 \rangle, \langle x_1, a_1, r_1, x_2 \rangle, \dots).$$

We differentiate between RL algorithms in two major ways. The first is whether the algorithm has control over its data: an algorithm is called *on-policy* if it controls its own actions from which it learns, and *off-policy* if it can learn from any data. Furthermore, we differentiate between *model-based* and *model-free* algorithms. Model-based algorithms learn the underlying MDP and solve it using value or policy iteration. Model-free algorithms only learn the value function, since, due to Bellman's theorem, that is all that is needed to act optimally.

10.1 Model-based

In model-based RL, we learn the MDP, *i.e.*, we estimate the transition probabilities $p(x' | x, a)$ and reward function $r(x, a)$ from the data,

$$\hat{p}(x' | x, a) = \frac{\text{count}(x' | x, a)}{\text{count}(x, a)}$$

$$\hat{r}(x, a) = \frac{1}{\text{count}(x, a)} \sum_{\substack{t=0 \\ x_t=x \\ a_t=a}}^{\infty} r_t$$

Then, we optimize the policy by value or policy iteration, based on the estimated MDP.

ϵ -greedy. At iteration t , pick random action with probability ϵ_t , or best action (according to internal MDP) with probability $1 - \epsilon_t$. Guaranteed to converge to optimal policy if the ϵ_t sequence satisfies the Robbins Monro conditions. The advantage of this method is that it is extremely simple



Figure 10.1. In reinforcement learning, an agent interacts with its environment. After playing an action a_t , it observes reward r_t and its new state x_{t+1} .

A sequence x_t satisfies the Robbins Monro conditions if

$$\sum_{t=0}^{\infty} x_t = \infty, \quad \sum_{t=0}^{\infty} x_t^2 < \infty.$$

E.g., $x_t = \frac{1}{t}$.

Algorithm	Classification		Space compl.
ϵ -greedy	On/off-policy	Model-based	$\mathcal{O}(\mathcal{A} \cdot \mathcal{X} ^2)$
R_{\max}	On/off-policy	Model-based	$\mathcal{O}(\mathcal{A} \cdot \mathcal{X} ^2)$
TD-learning	On-policy	Model-free	$\mathcal{O}(\mathcal{X})$
Q-learning	Off-policy	Model-free	$\mathcal{O}(\mathcal{A} \cdot \mathcal{X})$
Deep Q Network	Off-policy	Model-free	
REINFORCE	On-policy	Model-free	

Table 2. RL algorithms covered in this text with their types.

and has a clear interpretation w.r.t. the exploration-exploitation dilemma. The disadvantage is that it does not quickly eliminate clearly suboptimal actions. This is because it explores the state space in an uninformed manner. In other words, it explores while ignoring all past experience.

R_{\max} algorithm. R_{\max} solves the problem of ϵ -greedy by using the *Optimism in the face of uncertainty* principle. It assumes that any unexplored states are “fairy tale” states with high reward. More formally, if $r(x, a)$ is unknown, we set $\hat{r}(x, a) = R_{\max}$. Similarly, if $p(x' | x, a)$ is unknown, we set $\hat{p}(x' | x, a) = 1$ for some “fairy tale” state,

$$\begin{aligned}\hat{p}(x^* | x^*, a) &= 1 & \forall a \in \mathcal{A} \\ \hat{r}(x^*, a) &= R_{\max} & \forall a \in \mathcal{A}.\end{aligned}$$

This gives us an algorithm that has a bias toward exploring, but once it has explored a part of the state-action space, and observed it to be suboptimal, it can quickly eliminate it. Furthermore, the algorithm does not have to explicitly choose between exploration and exploitation, because it is done by assuming that the unexplored states are optimal.

10.2 Model-free

The problem with model-based RL is that it has high space requirements for storing the MDP, i.e., $\mathcal{O}(|\mathcal{A}| \cdot |\mathcal{X}|^2)$. Furthermore, it requires repeatedly solving the underlying MDP, which is expensive with policy or value iteration. In *model-free* RL, we estimate the value function directly, because that is all we need to act optimally, according to Bellman’s theorem. Thus, we also do not need to do any planning, eliminating much computational complexity.

Temporal difference-learning. TD-learning directly computes the value function. Recall the Bellman expectation equation,

$$V^\pi(x) = r(x, \pi(x)) + \gamma \sum_{x' \in \mathcal{X}} p(x' | x, \pi(x)) V^\pi(x').$$

Since we do not have access to r and p , we have to make a Monte Carlo estimate given a single data point $\langle x, a, r, x' \rangle$,

$$\approx r + \gamma V^\pi(x').$$

The idea is that we make this approximation repeatedly as the agent collects new data, which achieves the same effect as averaging over many data points. However, there is still a problem: V^π depends on the unknown V^π .

The key idea is to use a bootstrapping estimate of the value function. In other words, instead of the true value function V^π , we will use a running estimate \hat{V}^π . However, due to relying on a single sample, the value function will have a high variance, which is why we mix the new estimate with the previous one using a learning rate α_t ,

$$\hat{V}^\pi(x) \leftarrow (1 - \alpha_t)\hat{V}^\pi(x) + \alpha_t(r + \gamma\hat{V}^\pi(x')).$$

If the learning rate α_t satisfies the Robbins Monro conditions and all states are chosen infinitely often, \hat{V}^π is guaranteed to converge to the optimal value function V^π .

Note that, due to the Monte Carlo approximation w.r.t. transitions attained by following policy π , TD-learning is a fundamentally on-policy method. Further note that the space requirement of this algorithm is $\mathcal{O}(|\mathcal{X}|)$.

Q-learning. A generalization of TD-learning is Q-learning. Instead of directly learning the value function, which makes it inherently on-policy, it learns state-action values $Q(x, a)$. Then, we can compute $V^\pi(x) = \max_a Q(x, a)$. Like in TD-learning, we mix in the new estimate with the previous one according to learning rate α_t ,

$$\hat{Q}(x, a) \leftarrow (1 - \alpha_t)\hat{Q}(x, a) + \alpha_t \left(r + \gamma \max_{a' \in \mathcal{A}} \hat{Q}(x', a') \right).$$

The advantage of Q-learning is that it is off-policy, because the value is conditioned on the action. Thus, we can generate as much data as we need using a different algorithm, such as ϵ -greedy, and then estimate the Q-values from there.

Again, Q-learning is optimal if it satisfies the Robbins Monro conditions and all state-action pairs are chosen infinitely often. The space complexity of Q-learning is $\mathcal{O}(|\mathcal{A}| \cdot |\mathcal{X}|)$.

10.3 Model-free deep RL

The problem with all previously discussed methods is that they are only feasible in a small finite domain. In continuous domains, we would need an infinite amount of memory to store all values. Thus, we need to approximate this regression problem with function approximators, a.k.a. machine learning.

Let $V^\pi(x; \theta)$ be the function approximator that approximates the value function. Just like in TD-learning, we make the following Monte Carlo

In deep RL, the input could be anything that machine learning can process, e.g., video game frames with CNN or language with RNN.

estimation for a given data point $\langle x, a, r, x' \rangle$,

$$V^\pi(x) \approx r + \gamma V^\pi(x').$$

Then, we can define the loss function of our value function as the squared error from the true value function,

$$\ell(\theta; x, r, x') \doteq \frac{1}{2} (V^\pi(x; \theta) - V^\pi(x))^2,$$

which we estimate by using the Monte Carlo estimation,

$$\approx \frac{1}{2} \left(V^\pi(x; \theta) - \left(r + \gamma V^\pi(x'; \theta^{\text{old}}) \right) \right)^2.$$

The gradient of this loss is equal to the following,

$$\nabla_{V^\pi(x; \theta)} \ell(\theta; x, r, x') = V^\pi(x; \theta) - \left(r + \gamma V^\pi(x'; \theta^{\text{old}}) \right).$$

Using stochastic gradient descent, we then get the following update rule,

$$\begin{aligned} V^\pi(x; \theta) &\leftarrow V^\pi(x; \theta) - \alpha_t \left(V^\pi(x; \theta) - \left(r + \gamma V^\pi(x'; \theta^{\text{old}}) \right) \right) \\ &= (1 - \alpha_t) V^\pi(x; \theta) + \alpha_t \left(r + \gamma V^\pi(x'; \theta^{\text{old}}) \right), \end{aligned}$$

which is the same as the TD-learning update rule. Thus, the TD-learning update rule can be viewed as gradient descent on the squared loss!

Deep Q-network. The same result holds for the Q-learning update rule, where we do gradient descent on

$$\ell(\theta; x, a, x', r) = \frac{1}{2} \left(Q(x, a; \theta) - \left(r + \gamma \max_{a'} Q(x', a; \theta^{\text{old}}) \right) \right)^2.$$

This equation is called the *Bellman error*.

```

function DEEPQNETWORK( $\tau$ )
  initialize  $\theta$ 
  while not converged do
    pop  $\langle x, a, r, x' \rangle$  from  $\tau$ 
     $\theta \leftarrow \theta - \alpha_t \delta \nabla_{\theta} Q(x, a; \theta)$ 
    where  $\delta \doteq \left( Q(x, a; \theta) - \left( r + \gamma \max_{a'} Q(x', a; \theta^{\text{old}}) \right) \right)$ 
  end while
end function

```

Algorithm 5. Q-learning with function approximation.

However, this algorithm is quite slow to converge. To accelerate, we clone the network and maintain a constant “target” value across episodes, which we update once in a while.

Furthermore, deep Q networks suffer from “maximization bias”, which means that it tends to overestimate the actual Q value. This is caused by the max-operator used in the update rule, which also maximizes noise in the data. This is solved by the double deep Q network, where we use two target networks. We then take the minimum at each iteration of the two predictions by the two target networks.

Policy search methods. The problem with deep Q networks is that if the action-space is large or infinite, the $\max_{a'} Q(x', a; \theta)$ is no longer feasible. The solution to this is learning a parametrized policy $\pi(x; \theta)$, where the output is the action. In this case, we want to maximize the expected trajectory reward,

$$\begin{aligned} J(\theta) &= \mathbb{E}_{x_{0:T}, a_{0:T} \sim \pi_\theta} \left[\sum_{t=0}^T \gamma^t r(x_t, a_t) \right] \\ &= \mathbb{E}_{\tau \sim \pi_\theta} [r(\tau)]. \end{aligned}$$

Theorem 10.1. The following holds,

$$\nabla_\theta J(\theta) = \nabla_\theta \mathbb{E}_{\tau \sim \pi_\theta} [r(\tau)] = \mathbb{E}_{\tau \sim \pi_\theta} [r(\tau) \nabla_\theta \log \pi_\theta(\tau)].$$

Proof.

$$\begin{aligned} \nabla_\theta J(\theta) &= \nabla_\theta \int \pi_\theta(\tau) r(\tau) d\tau \\ &= \int \nabla_\theta \pi_\theta(\tau) r(\tau) d\tau \\ &= \int \pi_\theta(\tau) r(\tau) \nabla_\theta \log \pi_\theta(\tau) d\tau && \text{Chain rule} \\ &= \mathbb{E}_{\tau \sim \pi_\theta} [r(\tau) \nabla_\theta \log \pi_\theta(\tau)]. \end{aligned}$$

■

Using Theorem 10.1, we do not need to use the reparametrization trick to be able to compute gradients. We only need to compute $\nabla_\theta \log \pi_\theta(\tau)$. We can compute this gradient as follows,

$$\begin{aligned} \pi_\theta(\tau) &= p(x_0) \prod_{t=0}^T \pi_\theta(a_t | x_t) p(x_{t+1} | x_t, a_t) \\ \nabla_\theta \log \pi_\theta(\tau) &= \nabla_\theta \log p(x_0) + \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | x_t) + \sum_{t=0}^T \nabla_\theta \log p(x_{t+1} | x_t, a_t) \\ &= \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | x_t). \end{aligned}$$

So, to be able to compute gradients w.r.t. θ , we do not even need to know the underlying MDP! Putting this together, we get the following gradient,

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[r(\tau) \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | x_t) \right].$$

Even though these gradients are unbiased, they typically have large variance. We can reduce the variance by introducing *baselines*,

$$\mathbb{E}_{\tau \sim \pi_\theta} [r(\tau) \nabla_\theta \log \pi_\theta(\tau)] = \mathbb{E}_{\tau \sim \pi_\theta} [(r(\tau) - b(\tau)) \nabla_\theta \log \pi_\theta(\tau)],$$

```

1: function REINFORCE
2:   initialize  $\theta$ 
3:   repeat
4:     generate an episode  $\tau$  using  $\pi_\theta$ 
5:     for  $t = 1, \dots, T$  do
6:        $G_t \leftarrow r_t$ 
7:        $\theta \leftarrow \theta + \eta \gamma^t G_t \nabla_\theta \log \pi_\theta(a_t | x_t)$ 
8:     end for
9:   until done
10: end function

```

which have the same gradient. Thus, we are able to shift the reward up or down without influencing the gradient.

REINFORCE (Algorithm 6) sets its baseline adaptively to be the following,

$$b_t(\tau) = \sum_{t'=0}^{t-1} \gamma^{t'} r_{t'}.$$

Then, the gradient becomes the following,

$$\begin{aligned}
\nabla_\theta J(\theta) &= \mathbb{E}_{t \sim \pi_\theta} \left[\sum_{t=0}^T \left(\sum_{t'=0}^T \gamma^{t'} r_{t'} - \sum_{t'=0}^{t-1} \gamma^{t'} r_{t'} \right) \nabla_\theta \log \pi_\theta(a_t | x_t) \right] \\
&= \mathbb{E}_{t \sim \pi_\theta} \left[\sum_{t=0}^T \gamma^t \left(\sum_{t'=t}^T \gamma^{t'-t} r_{t'} \right) \nabla_\theta \log \pi_\theta(a_t | x_t) \right] \\
&= \mathbb{E}_{t \sim \pi_\theta} \left[\sum_{t=0}^T \gamma^t G_t \nabla_\theta \log \pi_\theta(a_t | x_t) \right],
\end{aligned}$$

where G_t is the reward to go following action a_t . REINFORCE is an on-policy algorithm, because it requires generating an episode for the data. This is necessary to be able to update its parameters correctly.

Actor-critic methods. We can reinterpret the REINFORCE gradient as follows,

$$\begin{aligned}
\nabla_\theta J(\theta) &= \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T \gamma^t G_t \nabla_\theta \log \pi_\theta(a_t | x_t) \right] \\
&= \sum_{t=0}^{\infty} \mathbb{E}_{\tau_{t:\infty}} [\gamma^t G_t \nabla_\theta \log \pi_\theta(a_t | x_t)] \\
&= \sum_{t=0}^{\infty} \mathbb{E}_{x_t, a_t} \left[\gamma^t \mathbb{E}_{\tau_{t:\infty}} \left[\sum_{t'=t}^{\infty} \gamma^{t'-t} r_{t'} \mid x_t, a_t \right] \nabla_\theta \log \pi_\theta(a_t | x_t) \right] \\
&= \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{\infty} \gamma^t Q^{\pi_\theta}(x_t, a_t) \nabla_\theta \log \pi_\theta(a_t | x_t) \right]
\end{aligned}$$

Algorithm 6. The REINFORCE algorithm, where the baseline at timestep t is set to be $\sum_{t'=0}^{t-1} \gamma^{t'} r_{t'}$.

$$r(\tau) - b_t = \sum_{t'=t}^T \gamma^{t'-t} r_{t'} \doteq G_t.$$

Intuitively, G_t is the reward to go following action a_t .



Figure 10.2. Illustration of an iteration of actor-critic methods.

Linearity of expectation

G_t depends on everything after t , while the other terms depend only on t .

The Q value is the value that we get after doing a_t at x_t .

Now, we can obtain the following,

$$\begin{aligned} &= \int \rho_{\theta}(x) \mathbb{E}_{a \sim \pi_{\theta}(x)} [Q(x, a) \nabla_{\theta} \log \pi_{\theta}(a | x)] dx \\ &= \mathbb{E}_{(x, a) \sim \pi_{\theta}} [Q(x, a) \nabla_{\theta} \log \pi_{\theta}(a | x)], \end{aligned}$$

This is an abuse of notation, because ρ_{θ} is an unnormalized probability distribution.

where $\rho_{\theta}(x) \doteq \sum_{t=0}^{\infty} \gamma^t p(x_t = x)$ is the unnormalized, discounted state occupancy measure.

This result naturally suggests plugging in approximations for $Q_{\theta}(x, a)$ for the action-value function. The idea of actor-critic networks is to parameterize an actor network that computes the policy and a critic network that computes the Q-value. They can then be used in each others' update equations,

$$\begin{aligned} \theta_{\pi} &\leftarrow \theta_{\pi} + \eta_t Q_{\theta}(x, a) \nabla_{\theta} \log \pi_{\theta}(a | x) \\ \theta_Q &\leftarrow \theta_Q - \eta_t (Q_{\theta}(x, a) - r - \gamma Q_{\theta}(x', \pi_{\theta}(x')) \nabla_{\theta} Q_{\theta}(x, a). \end{aligned}$$

Furthermore, we can introduce baselines by adding a value network,

$$\begin{aligned} \theta_{\pi} &\leftarrow \theta_{\pi} + \eta_t (Q_{\theta}(x, a) - V_{\theta}(x)) \nabla_{\theta} \log \pi_{\theta}(a | x) \\ &= \theta_{\pi} + \eta_t A(x, a) \nabla_{\theta} \log \pi_{\theta}(a | x), \end{aligned}$$

where $A(x, a) \doteq Q(x, a) - V(x)$ is the advantage function, which is positive if the chosen action is better than expected and negative if worse. Thus, intuitively, we increase the probability of the chosen action if better than expected, otherwise we decrease it. This model is called *advantage actor critic* (A2C) [Mnih et al., 2016].

All models discussed so far have been on-policy, which often causes sample inefficiency. Now, we want to move to off-policy methods. Recall that our initial motivation was that finding the maximum Q value was intractable if the action space was infinite. But, we could also replace the exact maximum by a parametrized policy,

$$J(\theta) = \sum_{(x, a, r, x') \in \mathcal{D}} \left(Q_{\theta}(x, a; \theta_Q) - \left(r + \gamma Q(x', \pi(x'; \theta_{\pi}); \theta_Q^{\text{old}}) \right) \right)^2,$$

where we jointly optimize over θ_Q and θ_{π} . We want to follow the greedy policy w.r.t. the Q function, *i.e.*, we want $\pi_{\theta} \approx \pi_Q = \arg\max_{a \in \mathcal{A}} Q(x, a; \theta_Q)$. The key idea is that if we use a “rich enough” parametrization of policies, selecting the greedy policy w.r.t. Q is equivalent to the following,

$$\theta_{\pi}^* = \arg\max_{\theta_{\pi}} \mathbb{E}_{x \sim \mu} [Q(x, \pi(x; \theta_{\pi}); \theta_Q)],$$

where $\mu(x) > 0$ is an *exploration distribution* over states with full support. If we then use differentiable approximations of Q and a differentiable deterministic policy π , we can use backpropagation to obtain gradients,

$$\nabla_{\theta_{\pi}} Q(x, \pi(x; \theta_{\pi}); \theta_Q) = \nabla_{\pi(x; \theta_{\pi})} Q(x, \theta_{\pi}) \nabla_{\theta_{\pi}} \pi(x; \theta_{\pi}).$$

However, policy gradient methods rely on randomized policies for exploration, but we have deterministic policies. To encourage exploration, we can inject additional Gaussian action noise to encourage exploration, akin to ϵ -greedy exploration. This is called the *deep deterministic policy gradients* (DDPG) algorithm [Lillicrap et al., 2015].

Twin delayed deep deterministic (TD3) further improves this by introducing a second critic network to address maximization bias [Fujimoto et al., 2018]. *Soft-actor critic* (SAC) further improves this by adding entropy regularization,

$$J_\lambda(\theta) = J(\theta) + \lambda H(\pi_\theta),$$

which encouraging exploration by giving preference to high-entropy actions [Haarnoja et al., 2018].

10.4 Model-based deep RL

So far, we have only discussed deep model-free methods. However, if we have an accurate model of the environment, we can use it for planning. The main benefit of this is that it dramatically reduces the sample complexity, compared to model-free techniques. In other words, we need much less data to find a good policy.

Planning in a known model. We assume a continuous state and action space with non-linear transitions, without constraints. Thus, this is quite a bit more complex than solving MDPs. We have a deterministic transition function f and a reward function r . The objective then becomes the following,

$$J_\infty(a_{0:\infty}) = \sum_{t=0}^{\infty} \gamma^t r(x_t, a_t).$$

such that $x_{t+1} = f(x_t, a_t)$ for every x_t .

However, we cannot plan over an infinite horizon. The key idea is to plan over a finite horizon H , carry out the first action, then replan with a horizon H . Thus, we first optimize over the following

$$J_H(a_{t:t+H-1}) = \sum_{t'=t}^{t+H-1} \gamma^{t'-t} r(x_{t'}, a_{t'}),$$

such that $x_{t+1} = f(x_t, a_t)$ for every x_t .

carry out action a_t and then replan. We can optimize this function using gradient methods (backpropagation through time) if the actions, rewards, and dynamics are differentiable. However, there are often many local minima, and vanishing/exploding gradients are a problem, because we do backpropagation through time.²⁹ Thus, we often use heuristic global optimization methods. The *random shooting* method generates m sets of random samples $a_{t:t+H-1}$, and then picks the sequence that optimizes the objective.

However, if we have sparse rewards, this will inherently not work, because we might not look far enough into the future. But, if we have

²⁹ This is the same problem that recurrent neural networks have.



Figure 10.3. Illustration of the effect of a finite horizon. Only after two steps, the agent could “see” the obstacle within its horizon.

access to a value function estimate, we can use that to see beyond the finite horizon, giving us the following object function,

$$J_H(a_{t:t+H-1}) = \sum_{t'=t:t+H-1} \gamma^{t'-t} r_{t'} + \gamma^H V(x_{t+H}),$$

where, intuitively, $V(x_{t+H})$ summarizes the remaining infinite timesteps. For $H = 1$, this is equal to the greedy policy w.r.t. V , but if we use larger H , this converges faster.

If the transition function is stochastic, rather than deterministic, we have to optimize over the expected performance,

$$J_H(a_{t:t+H-1}) = \mathbb{E}_{x_{t+1:t+H}} \left[\sum_{t'=t}^{t+H-1} \gamma^{t'-t} r_{t'} + \gamma^H V(x_{t+H}) \mid a_{t:t+H-1} \right].$$

However, the problem is that this expectation requires solving a high-dimensional integral. A common solution is to use *Monte Carlo trajectory sampling*. For this we use the reparametrization trick to obtain unbiased Monte Carlo estimates,

$$\hat{J}_H(a_{t:t+H-1}) = \frac{1}{m} \sum_{i=1}^m \sum_{t'=t}^{t+H-1} \gamma^{t'-t} r_{t'}(x_t(a_{t:t'-1}, \epsilon_{t:t'-1}^{(i)}), a_t) + \gamma^H V(x_{t+H}).$$

The state is a function of the previous actions $a_{t:t'-1}$ and noise $\epsilon_{t:t'-1}^{(i)}$ that is needed for the Gaussian reparameterization trick.

Instead of optimizing over the actions, we could also optimize over parametrized policies π_θ . This replaces expensive online planning by offline training of a policy that is fast to evaluate online. The objective then becomes the following,

$$J(\theta) = \mathbb{E}_{x_0 \sim \mu} \left[\sum_{t=0}^{H-1} \gamma^t r_t + \gamma^H Q(x_H, \pi_\theta(x_H)) \right].$$

For $H = 0$, this is equivalent to the DDPG objective.

Learning the model. Until now, we have assumed known f and r . In RL, these are unknown, so we have to learn them. The key insight is that due to the Markovian structure of the MDP, the observed transitions and rewards are conditionally independent. We can estimate them off-policy with standard supervised learning techniques from a dataset of

trajectories. We train a model with inputs (x_t, a_t) and the model outputs (r_t, x_{t+1}) . For continuous state spaces, this is essentially just a regression problem.

We could use an MAP estimate, but errors in the model estimate compound, which the planning algorithms exploit, which results in poor performance. This can easily be remedied by capturing uncertainty (Gaussian processes, Bayesian neural networks) in the estimated model and taking it into account in planning.

Exploration and exploitation. An algorithm that we can use to balance exploration and exploitation is Thompson sampling, which we have already seen before. We sample a model, plan a new policy according to it, roll out policy to collect more data, and finally update the posterior.

References

- Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International conference on machine learning*, pages 1587–1596. PMLR, 2018.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pages 1861–1870. PMLR, 2018.
- Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- Wesley J Maddox, Pavel Izmailov, Timur Garipov, Dmitry P Vetrov, and Andrew Gordon Wilson. A simple baseline for bayesian uncertainty in deep learning. *Advances in neural information processing systems*, 32, 2019.
- Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937. PMLR, 2016.