

Lecture 20: Gaussian Processes

What we did so far:

observe data: $[x_1, y_1], [x_2, y_2], \dots$

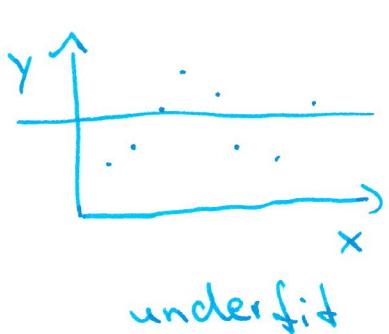
assume model, e.g. $y_i = \theta^T x_i + \epsilon_i$

estimate parameters: $p(\theta|D) \propto p(D|\theta) \cdot p(\theta)$

make predictions for new data inputs

$$p(y_* | X_*, x, y)$$

\Rightarrow A problem with this approach is coming up with the right model \Rightarrow might underfit or overfit



\Rightarrow It would be great if instead of estimating θ , we could estimate the model.

\Rightarrow We need a distribution over functions $f(x)$:

$$y_i = f(x_i) + \epsilon_i$$

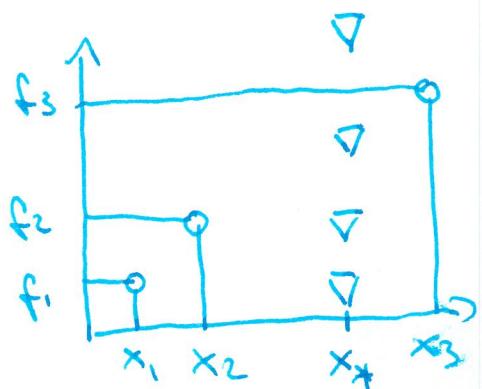
$$p(y_* | x_*, x, y) = \int p(y_* | f, x_*) \cdot p(f | x, y) df$$

\Rightarrow This looks rather difficult, but it is exactly what Gaussian Processes enable us to do!

Here is the trick: think about functions as lookup tables.

\Rightarrow We only need to be able to define a distribution over the function's values at a finite, but arbitrary set of points.

Let us look at an example:



which ∇ would you choose for f_* ?

\Rightarrow Intuitively we want our function to be smooth.

\Rightarrow points that are closer in x_i should be closer together in f_i .

$\Rightarrow x_1$ gives us information about x_2 , and a bit less information about x_3 .

In the Gaussian Process setting we model $p(f_1, f_2, f_3)$ as jointly Gaussian:

$$\begin{bmatrix} f_1 \\ f_2 \\ f_3 \end{bmatrix} \sim N \left(\underbrace{\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}}_{\mu}, \underbrace{\begin{bmatrix} 1 & 0.7 & 0.2 \\ 0.7 & 1 & 0.6 \\ 0.2 & 0.6 & 1 \end{bmatrix}}_K \right)$$

Note: These numbers are made up, it is about understanding the concept.

\Rightarrow Covariance captures similarity

$\Rightarrow \text{cov}(x_1, x_2) > \text{cov}(x_1, x_3)$

How can we construct the covariance matrix K ?

\Rightarrow We need some measure of similarity, and K needs to be positive definite.

$\Rightarrow K$ can be computed using standard covariance functions / kernels.

\Rightarrow One popular choice is the squared exponential:

$$K_{ij} = T^2 \cdot e^{-\frac{\|x_i - x_j\|^2}{2\lambda^2}}$$

\uparrow
not always
there

$\Rightarrow K_{ij} = 1$ if $x_i = x_j$

λ is a parameter that regulates the smoothness of our function / GP.

T regulates the height.

\Rightarrow Once again there is no free lunch. We have to decide for a kernel and tune the parameters.

\Rightarrow So far all x_i are given, and we can use our covariance function to construct K .

$\Rightarrow K$ is also called the Gram matrix.

How do we get a prediction f_* for a given new point x_* ?

Assumption: The new unseen data point comes from the same distribution as the training data:

$$f_* \sim N(0, K_{**})$$

$$K_{**} = \tau^2 e^{-\frac{\|x_* - x_*\|^2}{2\tau^2}} = 1$$

\Rightarrow If we define our problem like this, then the known f_i and f_* are independent, but we want smoothness, so they should be correlated.

$$\begin{bmatrix} f_1 \\ f_2 \\ f_3 \\ f_* \end{bmatrix} \sim N \left(\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} K_{11} & K_{12} & K_{13} & K_{1*} \\ K_{21} & K_{22} & K_{23} & K_{2*} \\ K_{31} & K_{32} & K_{33} & K_{3*} \\ K_{*1}, K_{*2}, K_{*3} & K_{*4} \end{bmatrix} \right)$$

\Rightarrow What we have now is a joint distribution

$$p(f_*, f_1, f_2, f_3 | X_*, X).$$

\Rightarrow What we want is the conditional

$$p(f_* | f, X_*, X).$$

\Rightarrow We remember HW5 and that the conditionals of a multivariate Gaussian are also Gaussian.

$$\Rightarrow p(f_* | f, X_*, X) \sim N(\mu_*, \Sigma_*)$$

$$\mu_* = \mu(X_*) + K_*^T K^{-1} (f - \mu(X))$$

$$\Sigma_* = K_{**} - K_*^T K^{-1} K_*$$

It is common to just set the mean of the joint distribution to zero, because the GP is flexible enough to model ~~not~~ arbitrary well. If you have a special case and really want to model the mean using a parametric function, please see Kevin P. Murphy, section 15.2.6.

So far we have looked at the noiseless case. What if we consider our noisy

$$y_i = f_i + \epsilon_i = f(x_i) + \epsilon_i ?$$

\Rightarrow The noise is independently added to each observation

\Rightarrow Noise only influences the diagonal of the covariance matrix

$$K_{ij} = K(x_i, x_j) + \sigma_e^2 \delta_{ij}$$

\hookrightarrow Kernel

\hookrightarrow Kronecker delta

$$\delta_{ij} = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j \end{cases}$$

Assuming the conventional zero mean for the joint distribution, we get the following conditionals:

$$p(f_*^* | X_*, X, Y) = N(\mu_*, \Sigma_*)$$

$$\mu_* = K_*^T K_Y^{-1} \cdot Y$$

$$\Sigma_* = K_{**} - K_*^T K_Y^{-1} K_*$$

$\Rightarrow K_Y$ is K plus the noise values for the diagonal.

A note about the matrix inverse:

Inverting a matrix can be tricky / numerically unstable. A more stable possibility is to use a Cholesky decomposition of the matrix K .

Cholesky decomposition:

$K = LL^T$ L is a lower triangular matrix with real and positive diagonal entries.

$$\begin{aligned}\mu_* &= K_*^T \cdot K_*^{-1} \cdot y \\ &= K_*^T \cdot (L \cdot L^T)^{-1} y \\ &= K_*^T \cdot L^{-T} \cdot \underbrace{L^{-1} \cdot y}_{\alpha}\end{aligned}$$

$m = L^{-1}y \Leftrightarrow Lm = y \Rightarrow$ solve linear system of equations to get m

$$\alpha = L^{-T} \cdot m \Leftrightarrow L^T \cdot \alpha = m$$

$$\Rightarrow \mu_* = K_*^T \cdot \alpha$$

Learning the Kernel parameter

marginal likelihood:

$$p(y|x, \theta) = \int p(y|f, x, \theta) \cdot p(f|x, \theta) df$$

$$p(f|x, \theta) = N(0, K)$$

$$p(y|f, \theta) = \prod_i N(f_i, \sigma_y^2)$$

$$\log(p(y|x, \theta)) = -\frac{1}{2} y' K_y^{-1} y - \frac{1}{2} \log |K_y| \\ - \frac{N}{2} \log (2\pi)$$

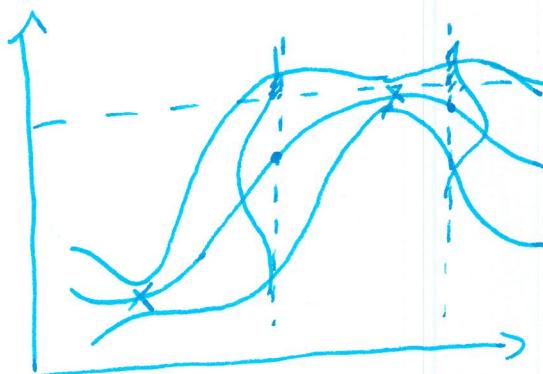
\Rightarrow We can use optimization methods and do maximum likelihood, or specify priors for the parameters θ and so for the posterior.

Active learning using GP

- Scenario: We want to find the maximum of a function, but querying points is expensive.
⇒ Need to balance exploration and exploitation

Example: drilling for oil
one armed bandit in a casino
winning a chess game

Probability of improvement:



$$f_{\max} = \mu^+ + \xi$$

↳ small value
↳ best value found
so far

PI looks at the area of the sample Gaussian that is above the f_{\max} line.

$$PI(x) = p(f(x) \geq \mu^+ + \xi)$$

\Rightarrow Offset ξ is necessary to deal with the case that we are looking at μ^+ again.

$$PI(x) = \phi\left(\frac{\mu(x) - \mu^+ - \xi}{\sigma(x)}\right) \quad \phi: \text{CDF of the Gaussian.}$$

PI is not often used in practice, but if you know the best you can get, it can be very powerful.

Expected Improvement

This criterion maximizes expected utility / minimizes expected cost.

Example:

$$p(x=\text{healthy} | \text{data}) = 0.9$$

$$p(x=\text{cancer} | \text{data}) = 0.1$$

	no treatment	treatment	$u(x, a)$
healthy	0	-30	
cancer	-100	-20	

$$EU(a) = \sum_x u(x, a) \cdot p(x | \text{data})$$

$$\begin{aligned} EU(a=\text{treatment}) &= u(\text{healthy, treatment}) \\ &\quad \cdot p(\text{healthy} | \text{data}) \\ &\quad + u(\text{cancer, treatment}) \\ &\quad \cdot p(\text{cancer} | \text{data}) \end{aligned}$$

\Rightarrow In this example the state is healthy or cancer, now the state is a function that we sample from our GP.

$$\underset{x}{\operatorname{argmin}} \int \|f_{n+1}(x) - f(x^*)\| \cdot p(f_{n+1} | D) df_{n+1}$$

↴ from ↴ true function,
 our GP maximum at
 x^*

\Rightarrow find the point that minimizes the expected distance to the true maximum over all

possible functions from our GP.

But we don't know the true function!

⇒ Use the best value we observed so far instead:

$$x_{n+1} = \underset{x}{\operatorname{argmax}} \text{IE} [\max \{0, f_{n+1}(x) - f^{\text{max}}\} | D_n]$$

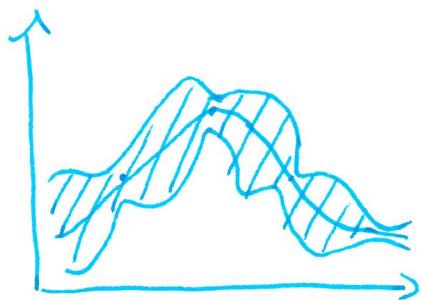
⇒ We want to do better than the f^{max} we have seen before.

It turns out that you can solve this expression analytically:

$$EI(x) = \begin{cases} (\mu(x) - \mu^* - \xi) \overset{\text{CDF}}{\uparrow} \phi(z) \cdot \sigma(x) \overset{\text{PDF}}{\uparrow} \phi(z) & \text{if } \sigma(x) > 0 \\ 0 & \text{if } \sigma(x) = 0 \end{cases}$$

$$z = \frac{\mu(x) - \mu^* - \xi}{\sigma(x)}$$

Thompson Sampling



If I sample a function from the GP, it has a high prob. of staying within the shaded area.

- ⇒ draw a function and choose the maximum of this function as the next point.
- ⇒ Good points create bottlenecks at good regions
⇒ exploitation.
- ⇒ Google analytics is using this.

Acknowledgement: This lecture is heavily based on Kevin P. Murphy, and on Nando de Freitas lectures in CPSC 540.

Nando is one of the best instructors I know and I highly recommend watching his YouTube videos.