

Gaussian Process Regression - Detailed Hand Calculations

1. Model Setup

Consider a simple 1D example with 3 training points:

- $X = [0, 1, 2]^T$
- $y = [0.5, 1.2, 0.8]^T$
- Test point: $x^* = 1.5$

Using squared exponential kernel with hyperparameters:

- $\kappa = 1.0$ (magnitude)
- $\ell = 1.0$ (lengthscale)
- $\sigma = 0.1$ (noise std dev)

2. Kernel Computations

2.1 Squared Exponential Kernel

$$k(x, x') = \kappa^2 \exp(-\frac{1}{2}(x-x')^2/\ell^2)$$

2.2 Training Covariance Matrix $K(X,X)$

$$K_{11} = k(0,0) = 1^2 \times \exp(0) = 1.000$$

$$K_{12} = k(0,1) = 1^2 \times \exp(-\frac{1}{2} \times 1^2/1^2) = \exp(-0.5) = 0.607$$

$$K_{13} = k(0,2) = 1^2 \times \exp(-\frac{1}{2} \times 4^2/1^2) = \exp(-2) = 0.135$$

$$K_{22} = k(1,1) = 1.000$$

$$K_{23} = k(1,2) = \exp(-0.5) = 0.607$$

$$K_{33} = k(2,2) = 1.000$$

$$K = \begin{bmatrix} 1.000 & 0.607 & 0.135 \\ 0.607 & 1.000 & 0.607 \\ 0.135 & 0.607 & 1.000 \end{bmatrix}$$

2.3 Adding Noise

$$C = K + \sigma^2 I = K + 0.01I$$

$$C = \begin{bmatrix} 1.010 & 0.607 & 0.135 \\ 0.607 & 1.010 & 0.607 \\ 0.135 & 0.607 & 1.010 \end{bmatrix}$$

2.4 Test-Train Covariance $K(x^*, X)$

$$k_1 = k(1.5, 0) = \exp(-\frac{1}{2} \times (1.5)^2 / 1^2) = \exp(-1.125) = 0.325$$

$$k_2 = k(1.5, 1) = \exp(-\frac{1}{2} \times (0.5)^2 / 1^2) = \exp(-0.125) = 0.882$$

$$k_3 = k(1.5, 2) = \exp(-\frac{1}{2} \times (0.5)^2 / 1^2) = \exp(-0.125) = 0.882$$

$$k^* = \begin{bmatrix} 0.325 & 0.882 & 0.882 \end{bmatrix}$$

3. Posterior Mean Calculation

3.1 Solve $C^{-1}y = \alpha$

$$C \times \alpha = y$$

$$\begin{bmatrix} 1.010 & 0.607 & 0.135 \end{bmatrix} \begin{bmatrix} \alpha_1 \end{bmatrix} = \begin{bmatrix} 0.5 \end{bmatrix}$$

$$\begin{bmatrix} 0.607 & 1.010 & 0.607 \end{bmatrix} \begin{bmatrix} \alpha_2 \end{bmatrix} = \begin{bmatrix} 1.2 \end{bmatrix}$$

$$\begin{bmatrix} 0.135 & 0.607 & 1.010 \end{bmatrix} \begin{bmatrix} \alpha_3 \end{bmatrix} = \begin{bmatrix} 0.8 \end{bmatrix}$$

Using matrix inversion (or linear system solver):

$$\alpha \approx \begin{bmatrix} 0.127 \end{bmatrix}$$

$$\dots \begin{bmatrix} 0.874 \end{bmatrix}$$

$$\dots \begin{bmatrix} 0.332 \end{bmatrix}$$

3.2 Compute Posterior Mean

$$\mu^* = k^{*T} \times \alpha$$

$$\dots = \begin{bmatrix} 0.325 & 0.882 & 0.882 \end{bmatrix} \times \begin{bmatrix} 0.127 \end{bmatrix}$$

$$\dots \begin{bmatrix} 0.874 \end{bmatrix}$$

$$\dots \begin{bmatrix} 0.332 \end{bmatrix}$$

$$\dots = 0.325 \times 0.127 + 0.882 \times 0.874 + 0.882 \times 0.332$$

$$= 0.041 + 0.771 + 0.293$$

$$\dots = 1.105$$

4. Posterior Variance Calculation

4.1 Prior Variance at Test Point

$$k^{**} = k(1.5, 1.5) = 1.000$$

4.2 Solve $C^{-1}k^* = v$

$$C \times v = k^*$$

Solution:

$$v \approx \begin{bmatrix} 0.007 \\ 0.744 \\ 0.241 \end{bmatrix}$$

4.3 Compute Posterior Variance

$$\begin{aligned} \sigma^{2*} &= k^{**} - k^{*T} \times v \\ &= 1.000 - [0.325 \quad 0.882 \quad 0.882] \times \begin{bmatrix} 0.007 \\ 0.744 \\ 0.241 \end{bmatrix} \\ &= 1.000 - (0.002 + 0.656 + 0.213) \\ &= 1.000 - 0.871 \\ &= 0.129 \end{aligned}$$

4.4 Standard Deviation

$$\sigma^* = \sqrt{0.129} = 0.359$$

5. Predictive Distribution

For observation y^* at $x^* = 1.5$:

$$\begin{aligned} p(y^*|X,y) &= N(1.105, 0.129 + 0.01) \\ &= N(1.105, 0.139) \end{aligned}$$

95% confidence interval:

$$\begin{aligned} [\mu^* - 2\sigma^*, \mu^* + 2\sigma^*] &= [1.105 - 2 \times 0.373, 1.105 + 2 \times 0.373] \\ &= [0.359, 1.851] \end{aligned}$$

6. Marginal Likelihood

6.1 Cholesky Decomposition

$C = LL^T$ where L is lower triangular

6.2 Log Marginal Likelihood

$$\log p(y|X, \theta) = -\frac{1}{2}y^T C^{-1}y - \frac{1}{2}\log|C| - \frac{1}{2}n \log(2\pi)$$

Components:

1. Data fit term: $-\frac{1}{2}y^T C^{-1}y$
2. Complexity penalty: $-\frac{1}{2}\log|C| = -\sum \log(L_{ii})$
3. Normalization: $-\frac{1}{2}n \log(2\pi)$

7. Kernel Properties

Squared Exponential Properties:

- Infinitely differentiable (very smooth)
- Universal approximator
- Stationary and isotropic
- Characterized by lengthscale ℓ

Effect of Hyperparameters:

- **κ (magnitude)**: Controls function variance
 - Large $\kappa \rightarrow$ large function variations
 - Small $\kappa \rightarrow$ small function variations
- **ℓ (lengthscale)**: Controls smoothness
 - Large $\ell \rightarrow$ smooth, slowly varying functions
 - Small $\ell \rightarrow$ rough, quickly varying functions
- **σ (noise)**: Controls uncertainty
 - Large $\sigma \rightarrow$ high observation noise
 - Small $\sigma \rightarrow$ low observation noise

8. Computational Complexity

- Training: $O(n^3)$ due to matrix inversion
- Prediction mean: $O(n^2)$ per test point
- Prediction variance: $O(n^2)$ per test point
- Storage: $O(n^2)$ for kernel matrix

9. Practical Considerations

1. Numerical Stability:

- Add jitter (small diagonal term) to K
- Use Cholesky decomposition
- Avoid explicit matrix inversion

2. Hyperparameter Optimization:

- Maximize log marginal likelihood
- Use gradient-based optimization
- Consider multiple random restarts

3. Scalability:

- For large datasets, use sparse GPs
- Inducing points approximation
- Local approximations