

STAT 518

Update Student Presentation

Wen Wei Loh

April 30, 2013

What's the paper about again?

- Gaussian Processes (GP)
 - Focus on **covariance functions** between outcomes
 - how two outcomes are correlated, based on values of their predictors
 - letting the data speak for itself
- Prediction
 - Objective is to obtain a **predictive** distribution for the outcome of a future observation
- Radford M. Neal [1999]

Yes, there was some math ...

- Simple Linear Regression Example
 - Observed data : $(x^{(1)}, t^{(1)}) , \dots , (x^{(n)}, t^{(n)})$
 - $t^{(i)}$ = outcome (target) for case i
 - $x^{(i)}$ = fixed input (predictor) for case i
- univariate for this example, but may be p -dimensional
- Fit a familiar model:

$$t^{(i)} = \alpha + x^{(i)}\beta + \epsilon^{(i)}, \quad \epsilon^{(i)} \underset{iid}{\sim} N(0, \sigma_{\epsilon}^2)$$

- Put independent priors on unknown parameters α and β :

$$\alpha \sim N(0, \sigma_{\alpha}^2), \quad \beta \underset{iid}{\sim} N(0, \sigma_{\beta}^2)$$

Yes, there was some math ...

⇒ Prior joint **multivariate Gaussian** distribution for t :

$$t = \left[t^{(1)}, \dots, t^{(n)} \right]^T \sim N_n(\mathbf{0}, \mathbf{C})$$

Covariance matrix $\mathbf{C} = \left\{ \text{Cov} \left[t^{(i)}, t^{(j)} \right] \mid i, j \in [1, n] \right\}$

$$\text{Cov} \left[t^{(i)}, t^{(j)} \right] = \text{E} \left[\left(\alpha + x^{(i)} \beta + \epsilon^{(i)} \right) \left(\alpha + x^{(j)} \beta + \epsilon^{(j)} \right) \right]$$

$$= \sigma_\alpha^2 + x^{(i)} x^{(j)} \sigma_\beta^2 + \delta_{ij} \sigma_\epsilon^2,$$

$\delta_{ij} = 1$ if $i = j$ and 0 otherwise i.e. Kronecker delta

So why the fuss over the covariance function?

$$C_{ij} = \text{Cov} [t^{(i)}, t^{(j)}] = \underbrace{\sigma_\alpha^2}_{\text{Constant}} + \underbrace{x^{(i)} x^{(j)} \sigma_\beta^2}_{\text{Linear}} + \underbrace{\delta_{ij} \sigma_\epsilon^2}_{\text{Noise}}$$

- Covariance between the outcomes is fully described by the **relationship between** the predictors

e.g. Linear function of predictors: restrictive?

Suppose the covariance depends on **differences** between predictor values, rather than just the values themselves?

⇒ Are there other ways to describe the covariance function?

So why the fuss over the covariance function?

- How about an exponential term instead?

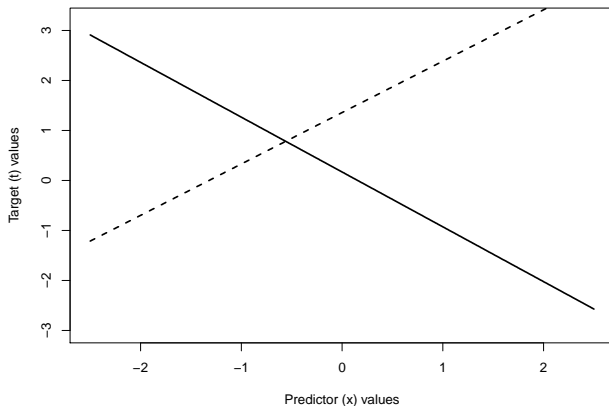
$$C_{ij} = \text{Cov} [t^{(i)}, t^{(j)}] = \underbrace{\eta^2 \exp \left(-\rho^2 (x^{(i)} - x^{(j)})^2 \right)}_{\text{Exponential}} + \underbrace{\delta_{ij} \sigma_\epsilon^2}_{\text{Noise}}$$

- Observations with predictor values that are “far apart” have much smaller covariances
- η = magnitude
- ρ = scale

\Rightarrow Gaussian kernel: $\exp \left(-\frac{(x^{(i)} - x^{(j)})^2}{1/\rho^2} \right)$

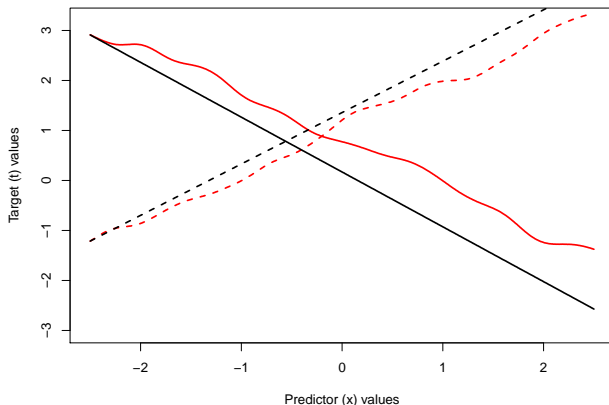
Show me some graphs instead

$$\text{Cov}[t^{(i)}, t^{(j)}] = 1 + x^{(i)}x^{(j)}$$



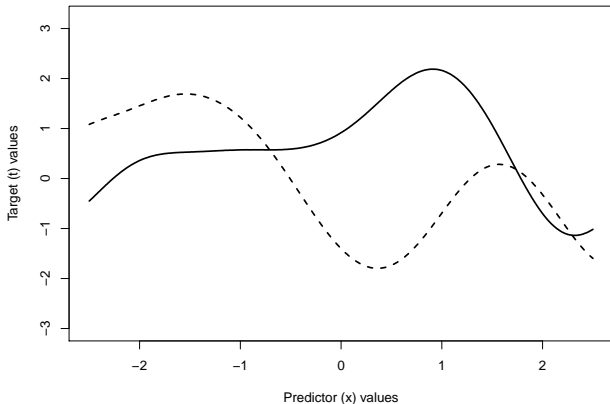
Show me some graphs instead

$$\text{Cov}[t^{(i)}, t^{(j)}] = 1 + x^{(i)}x^{(j)} + 0.1^2 \exp\left(-3^2 (x^{(i)} - x^{(j)})^2\right)$$
$$\text{Cov}[t^{(i)}, t^{(j)}] = 1 + x^{(i)}x^{(j)}$$



Show me some graphs instead

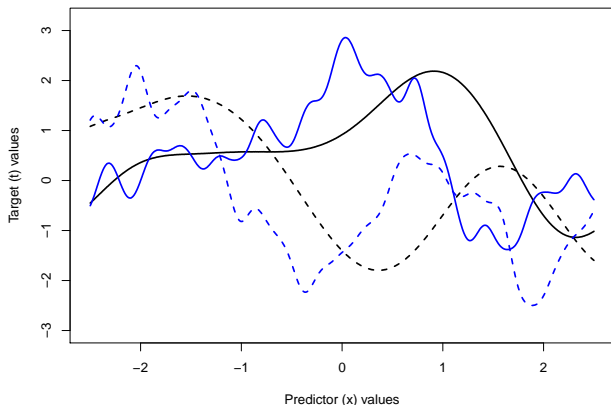
$$\text{Cov} [t^{(i)}, t^{(j)}] = \exp \left(- (x^{(i)} - x^{(j)})^2 \right)$$



Show me some graphs instead

$$\text{Cov}[t^{(i)}, t^{(j)}] = \exp\left(-\left(x^{(i)} - x^{(j)}\right)^2\right) + \mathbf{0.5^2} \exp\left(-\mathbf{5^2} \left(x^{(i)} - x^{(j)}\right)^2\right)$$

$$\text{Cov}[t^{(i)}, t^{(j)}] = \exp\left(-\left(x^{(i)} - x^{(j)}\right)^2\right)$$



Enough with the covariance functions ...

- May be constructed with various components e.g. constant, linear, exponential.
- Components of the covariance function may reflect different plausible features of the underlying structure
- Valid covariance function must always result in a **positive definite** covariance matrix for the targets.
- Different forms of the **covariance function** $\text{Cov} [t^{(i)}, t^{(j)}]$ define infinitely many flexible regression models

Step back to prediction

- Targets **still** have a **multivariate Gaussian** distribution

$$t = [t^{(1)}, \dots, t^{(n)}]^T \sim N_n(\mathbf{0}, \mathbf{C})$$


- $\mathbf{C}_{ij} = \text{Cov} [t^{(i)}, t^{(j)}]$
- Given the inputs for a new case $x^{(n+1)}$, the predictive distribution for the new outcome $t^{(n+1)}$ is **Gaussian**:

$$\begin{aligned} \mathbb{E} [t^{(n+1)} \mid t^{(1)}, \dots, t^{(n)}] &= k^T \mathbf{C}^{-1} t \\ \text{var} [t^{(n+1)} \mid t^{(1)}, \dots, t^{(n)}] &= V - k^T \mathbf{C}^{-1} k \end{aligned}$$

- $k = \left(\text{Cov} [t^{(n+1)}, t^{(1)}], \dots, \text{Cov} [t^{(n+1)}, t^{(n)}] \right)^T$
- $V = \text{Cov} [t^{(n+1)}, t^{(n+1)}] = \text{prior var} [t^{(n+1)}]$

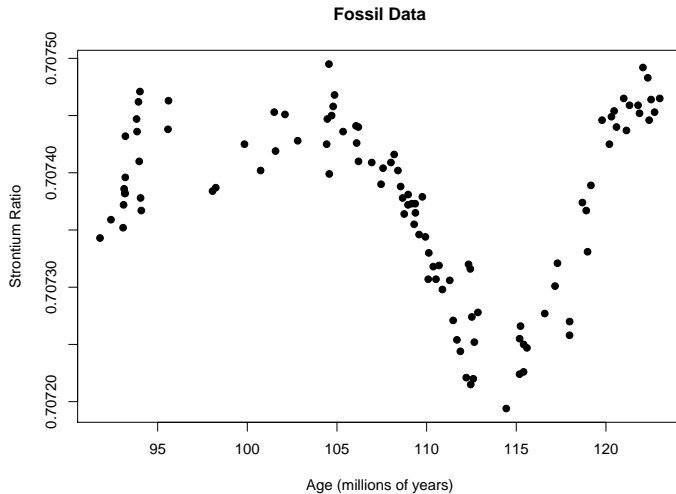
Finally ... some data

- 106 observations of fossil shells
- Age (in millions of years) and Ratios of strontium isotopes
- Previous example from STAT 527
- Data from SemiPar package in R [Ruppert et al., 2003], who got it from Bralower et al. [1997]



⁸⁴ Sr 83.913426 0.56% Stable	⁸⁶ Sr 85.909265 9.86% Stable	⁸⁷ Sr 86.908882 7.00% Stable	⁸⁸ Sr 87.905617 82.58% Stable
--	--	--	---

An “old” dataset



How to fit a Gaussian Process Regression Model

- Assume the prior covariance function:

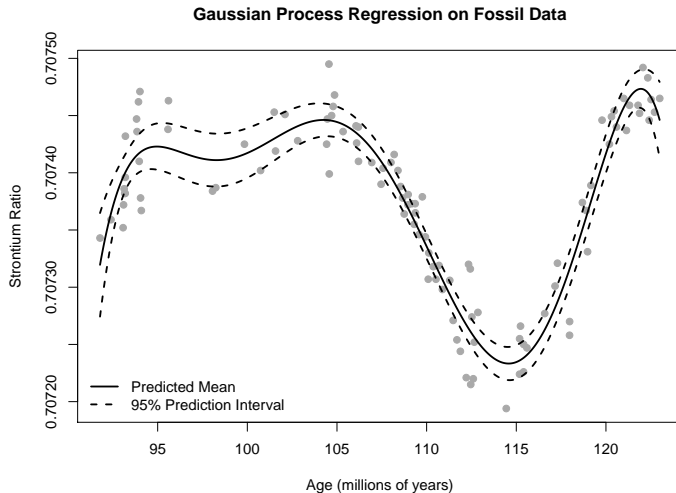
$$\mathbf{C}_{ij} = \text{Cov} [t^{(i)}, t^{(j)}] = \eta^2 \exp \left(-\rho^2 (x^{(i)} - x^{(j)})^2 \right) + \delta_{ij} \sigma_\epsilon^2$$

- Based on the multivariate Gaussian distribution of the outcomes, the log-likelihood is:

$$\log p(t \mid x, \eta, \rho) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log |\mathbf{C}| - \frac{1}{2} t^T \mathbf{C}^{-1} t$$

- Find maximum likelihood estimates $\hat{\eta}, \hat{\rho}$
- σ_ϵ^2 assumed to be fixed at 10^{-9}

Gaussian Process Regression on Fossil data



What's Next?

$$C_{ij} = \text{Cov} [t^{(i)}, t^{(j)}] = \eta^2 \exp \left(-\rho^2 (x^{(i)} - x^{(j)})^2 \right) + \delta_{ij} \sigma_\epsilon^2$$

- How to implement a Bayesian approach?
 - Integrate out parameters to get “parameter-free” marginal distribution $p(t^{(n+1)} | t^{(1)}, \dots, t^{(n)})$ for prediction?
- More to come ...
 - Compare against other semi-parametric procedures
 - p -dimension covariate regression example ($p > 1$)
 - How to implement a three-way classification / discrimination procedure?

References

- T J Bralower, P D Fullagar, C K Paull, G S Dwyer, and R M Leckie. Mid-Cretaceous strontium-isotope stratigraphy of deep-sea sections. *Geological Society of America Bulletin*, 109:1421–1442, 1997. doi: 10.1130/0016-7606(1997)109<1421.
- R M Neal. Regression and classification using Gaussian process priors (with discussion). *Bayesian Statistics*, 6: 475–501, 1999.
- David Ruppert, Matthew P Wand, and Raymond J Carroll. *Semiparametric regression*, volume 12. Cambridge University Press, 2003. ISBN 0521785162.