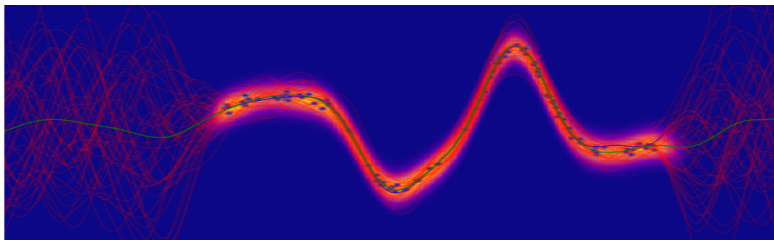


Stochastische Regressionsmodelle für heterogene Messnetzwerke

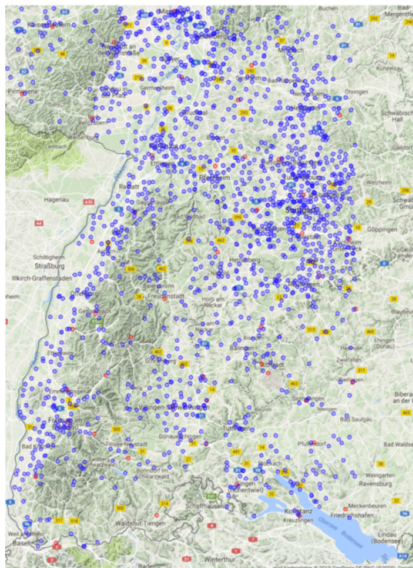
Gaußprozesse, Kernel, Hyperparameteroptimierung.

Dr. rer. nat. Johannes Riesterer - KIT/TECO



Automated Quality Assessment of (Citizen) Weather Stations; Julian Bruns, Johannes Riesterer, Bowen Wang, Till Riedel, Michael Beigl ; GI-Forum (2018).

Motivation



Gegeben

Gegeben sehr grosse Stichproben (Big Data) $D = \{(y_i, x_i)_i\}$ der Zufallsvariablen $Y : \mathbb{R}^n \rightarrow \mathbb{R}$, $X : \mathbb{R}^n \rightarrow \mathbb{R}^m$ mit unbekanntem Verteilungen.

Gesucht

Funktion $f : \mathbb{R}^m \rightarrow \mathbb{R}$ mit $f = \operatorname{argmin}_h L(Y, h(X))$, wobei L eine sogenannte Fehler-Funktion ist.

Bedingter Erwartungswert

$$L(Y, h(X)) := \int (Y - h(X))^2 d\mu \Rightarrow f = \mathbb{E}(Y|X)$$

Conclusio

Bestimme "posterior predictive distribution" $p(y|x, D)$

Algorithmisches Problem

$$Y = f(X) =: f_X$$

$$f = \operatorname{argmin}_{\{f_\theta \mid \theta \in \Theta\}} L(Y, f_\theta(X))$$

$$L(Y, h(X)) \doteq \sum_D (y - h(x))^2$$

$$D = D_{\text{Train}} \cup D_{\text{Test}} \cup D_{\text{Validation}}$$

Posterior predictive distribution

$$\begin{aligned} p(y^*|x^*, D) &= p(f_{x^*}|x^*, D) = \int_{\Omega} p(f_{y^*}, \omega|x^*, D) d\omega = \\ &= \int_{\Omega} \underbrace{p(f_{x^*}|\omega, x^*, D)}_{=p(f_{x^*}|\omega, x^*)} \cdot p(\omega|D) d\omega \end{aligned}$$

Satz von Bayes

$$\begin{aligned} \text{posterior} &= \frac{\text{likelihood} \cdot \text{prior}}{\text{marginal likelihood}} \\ p(\omega|D) &= \frac{p(D|\omega) \cdot p(\omega)}{p(D)} \end{aligned}$$

Conclusio

Wähle problemspezifische likelihood und prior Verteilungen $p(D|\omega)$ und $p(\omega)$.

Konjugierte Paare

Selten analytisch lösbar - wichtige Ausnahme konjugierte Paare.

Rejection sampling/Markov-Chain-Monte-Carlo

Integrale lösen durch sampling:

$$\sum_{\theta_i} \frac{g(\theta_i)}{p(\theta_i)} \rightarrow \int_{\Theta} g d\theta$$

Normalisierung

Mit der Marginalisierung ist

$p(D) = \int_{\omega'} p(D, \omega') d\omega' = \int_{\omega'} p(D|\omega') \cdot p(\omega') d\omega'$ nun berechenbar.

Modell

$Y = f(X) = X^t \cdot \omega + \epsilon$ mit Gewichten $\omega \sim \mathcal{N}(0, \Sigma)$,
 $\epsilon \sim \mathcal{N}(0, \sigma^2)$ und X konstant.

Likelihood

$$(y_i | x_i, \omega) \sim \mathcal{N}(\omega^t \cdot x_i, \sigma^2)$$

$$\Leftrightarrow p(y_i | x_i, \omega) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(y_i - \omega^t x_i)^2}{2\sigma}\right)$$

$$y_i \text{ u.i.v.} \Rightarrow p(Y|X, \omega) = \prod_i \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(y_i - \omega^t x_i)^2}{2\sigma}\right)$$

Posterior distribution

$$p(\omega|S) = \frac{(\prod_i p(y_i|x_i, \omega)) \cdot p(\omega)}{\int_{\omega'} (\prod_i p(y_i|x_i, \omega')) \cdot p(\omega') d\omega'}$$

$$\Rightarrow \omega|S \sim \mathcal{N}\left(\frac{1}{\sigma^2} A^{-1} X^t Y, A^{-1}\right)$$

mit $A = \frac{1}{\sigma^2} X^t X + \Sigma^{-1}$.

Posterior predictive distribution

$$y^*|x^*, D \sim \mathcal{N}\left(\frac{1}{\sigma^2} (x^*)^t A^{-1} X^t Y, (x^*)^t A^{-1} x^*\right)$$

mit $A = \frac{1}{\sigma^2} X^t X + \Sigma^{-1}$

Vorhersage

Für Feature x^* wird die Vorhersage durch

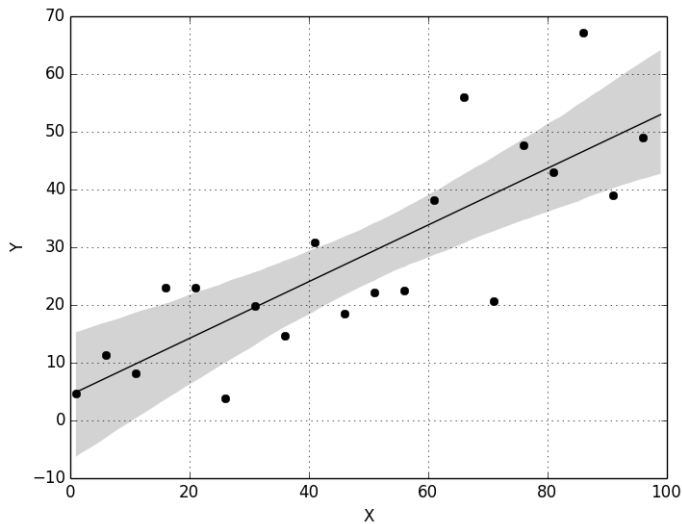
$$\mathbb{E}((y^*|x^*, D)) = \frac{1}{\sigma^2}(x^*)^t A^{-1} X^t Y$$

definiert. Die Varianz

$$\mathbb{V}((y^*|x^*, D)) = (x^*)^t A^{-1} x^*$$

dient als Mass der Güte der Vorhersage.

Lineare Bayessche Regression



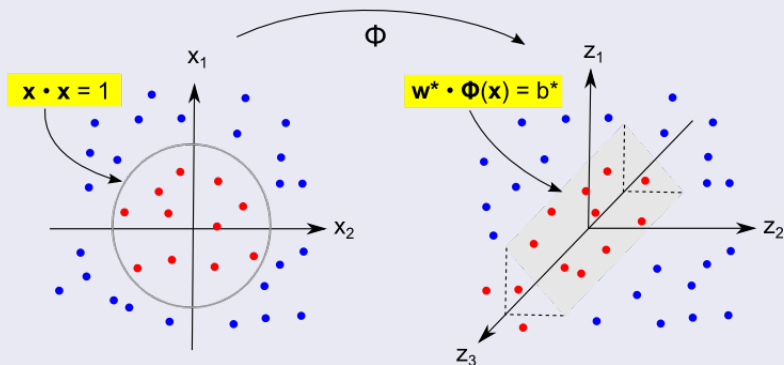
Bayessche Regression

Kernel Trick

$$f(\mathbf{X}) = \phi(\mathbf{X})^t \cdot \omega \text{ mit } \phi : \mathbb{R}^n \rightarrow \mathbb{R}^N$$

Beispiel

$$\Phi(\mathbf{x}) = \Phi(x_1, x_2) = (x_1^2, x_2^2, \sqrt{2}x_1x_2) = (z_1, z_2, z_3) = \mathbf{z}$$



Bayessche Regression

Analoge Rechnung zu linearer Bayesscher Regression:

$$f_{x^*} | x^*, D \sim \mathcal{N} \left(\frac{1}{\sigma^2} \underbrace{\phi(x^*)^t A^{-1} \phi(x)}_{:=k(x^*, x)} Y, \phi(x^*)^t A^{-1} \phi(x^*) \right)$$

mit $A = \frac{1}{\sigma^2} \phi(x)^t \phi(x) + \Sigma^{-1}$

Kernel

Funktionen der Form

$$k(x^*, x) = \psi(x^*)^t \psi(x)$$

werden Kernel oder Kovarianzfunktionen genannt.

Mercers Theorem

Ist $k(s, t) = k(t, s)$ und $\int_{X \times X} k(s, t) f(s) f(t) > 0$ für alle $f \in L^2(X)$, dann ist

$$k(s, t) = \sum_{i=1}^{\infty} \lambda_i \phi_i(s) \phi_i(t)$$

wobei ϕ_i die Eigenfunktionen zum Eigenwert λ_i des linearen Operators $T_k f := \int_X k(\cdot, t) f(t)$ sind.

Stochastischer Prozess

Ein stochastischer Prozess ist eine indizierte Menge von Zufallsvariablen $\{f_x \mid x \in \mathcal{X}\}$.

Gauß-Prozess

Wir bezeichnen einen stochastischen Prozess als Gauß-Prozess $f_x \sim \mathcal{GP}(m(x), k(x, x'))$, falls

$$\begin{pmatrix} f_{x_1} \\ \vdots \\ f_{x_n} \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} m(x_1) \\ \vdots \\ m(x_n) \end{pmatrix}, \begin{pmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{pmatrix} \right)$$

für jede endliche Teilmenge $X = (x_1, \dots, x_n) \in \mathcal{X}$. Man nennt $k(x, x')$ Kovarianz-Funktion oder auch Kernel. Zulässig sind nur Funktionen, bei denen die Matrix positiv definit und symmetrisch ist für jede endliche Teilmenge X .

Beispiel

$f_x = \phi(x)^t \cdot \omega$ mit $\omega \sim \mathcal{N}(0, \Sigma)$ ist ein $\mathcal{GP}(m(x), k(x, x'))$ mit

$$\mathbb{E}(f_x) = \phi(x)^t \mathbb{E}(\omega) = \underbrace{0}_{=: m(x)}$$

$$\mathbb{E}(f_x f_{x'}) = \phi(x)^t \mathbb{E}(\omega \omega^t) \phi(x') = \underbrace{\phi(x)^t \Sigma \phi(x')}_{:= k(x, x')}$$

Prior distribution

Sei $f \sim \mathcal{GP}(0, k(x, x'))$ ein Gauß-Prozess. Angenommen man kennt $f = (f_{x_1} \dots f_{x_n})$ an den Punkten $X = (x_1, \dots, x_n)$ und möchte $f^* = (f_{x_1^*} \dots f_{x_n^*})$ an den Punkten $X^* = (x_1^*, \dots, x_n^*)$ vorhersagen. Aus der GP-Eigenschaft folgt

$$\begin{pmatrix} f_X \\ f_{X^*} \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} K(X, X) & K(X, X^*) \\ K(X^*, X) & K(X^*, X^*) \end{pmatrix} \right)$$

Posterior predictive distribution

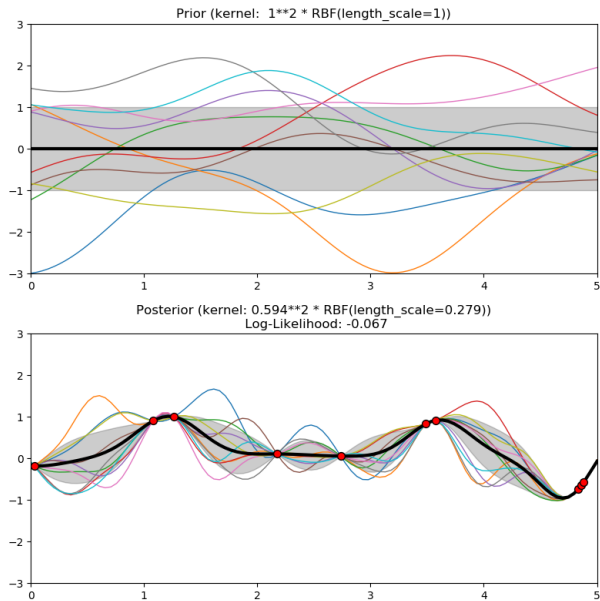
Mit den Rechenregeln für multivariate Normalverteilungen folgt:

$$f_{X^*} | X^*, f_X, X \sim \mathcal{N}(\mu, \Sigma)$$

$$\mu := K(X^*, X)K(X, X)^{-1} \cdot f$$

$$\Sigma := K(X^*, X^*) - K(X^*, X)K(X, X)^{-1}K(X, X^*)$$

Gauß-Prozess-Regression



Vorhersage

Für Prior \tilde{f} wird die Vorhersage durch

$$R(\tilde{f}) := \mathbb{E}((\tilde{f}|\tilde{X}, X, f)) = K(\tilde{X}, X)K(X, X)^{-1} \cdot f$$

definiert. Die Varianz

$$\mathbb{V}((\tilde{f}|\tilde{X}, X, f)) = K(\tilde{X}, \tilde{X}) - K(\tilde{X}, X)K(X, X)^{-1}K(X, \tilde{X})$$

dient als Mass der Güte der Vorhersage.

Kernel

Kernel	Funktion
konstant	σ_0^2
linear	$\sum_{d=1}^D \sigma_d^2 x_d x'_d$
polynomial	$(x \cdot x' + \sigma_0^2)^p$
squared exponential	$\exp\left(-\frac{r^2}{2l^2}\right)$
Matérn	$\frac{1}{2^{\nu-1} \Gamma(\nu)} \left(\frac{\sqrt{2\nu} r}{l}\right)^{\nu} K_{\nu} \left(\frac{\sqrt{2\nu} r}{l}\right)$
exponentiell	$\exp\left(-\frac{r}{l}\right)$
γ -exponentiell	$\exp\left(-\left(\frac{r}{l}\right)^{\gamma}\right)$
rational quadratisch	$\left(1 + \frac{r^2}{2\alpha l^2}\right)^{-\alpha}$

Qualitätsmodell

Gegeben Menge von Sensoren $\{(p_i, y_i)\}$ mit Position p_i und Messwert y_i unbekannter Qualität. Weise jedem Sensor einen Parameter $q_i \in (0, \infty)$ zu.

Kernel

$$K((p_i, q_i), (p_j, q_j)) := K_{\text{Matern}}(p_i, p_j) + K_q(q_i, q_j)$$

$$K_q(q_i, q_j) := \begin{cases} \frac{1}{q_i^2} & i=j \\ 0 & \text{else} \end{cases}$$

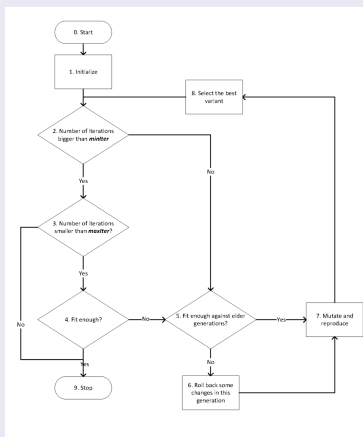
Modell

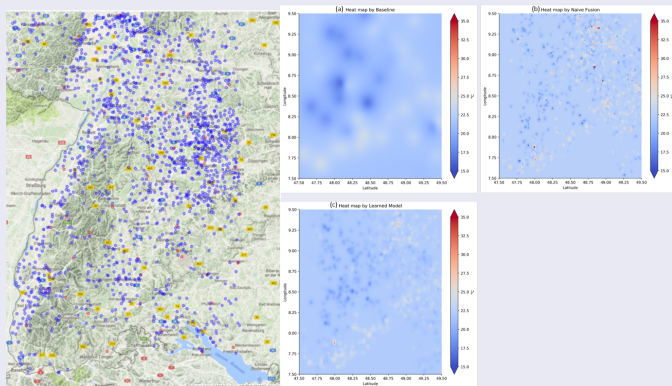
Erhalten damit Regressionmodell $GP(Q)$ mit Hyperparametern $Q = \{(q_i)\}$

Optimierung

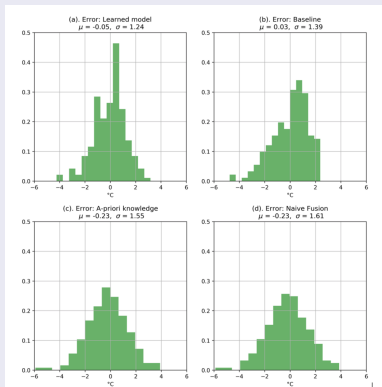
Optimiere $GP(Q)$ bzgl Q . Zum Beispiel mit Hilfe eines genetischen Algorithmus.

Ablaufdiagramm





Hypothesentest



Signifikant zum Signifikanzniveau $\alpha = 0.05$