

Oh!! Gaussian Process building

So... basic formulae:

$$\vec{y} = f(\vec{x}) + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I), \quad , y = \text{measurements (noisy)}$$

$$\Rightarrow p(y | f(x)) \sim N(f(x), \sigma^2) \quad + \text{then } f(x)$$

then, $p(f | y, x)$ = posterior
 say ~~$f(x) \sim N(\mu_f, \Sigma_f)$~~
 by prior

Standard
Linear
Model

$$y = f(x) + \varepsilon = \vec{x}^T w + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I)$$

$$w \sim N(\mu_w, \Sigma_w)$$

$$\vec{x} = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \end{pmatrix} \rightarrow \vec{x} = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{D1} & x_{D2} & \dots & x_{Dn} \end{pmatrix} \quad D$$

$$\sim A(\vec{x}, \sigma^2)$$

$$p(y | X, w, \sigma^2) = N(\vec{x}^T w, \sigma^2 I)$$

add mean & variance of Gaussians
 (ε only varying parameter, mean=0)

↳ likelihood.

$$\text{prior: } p(w) = N(0, \Sigma_w)$$

$$\text{posterior: } p(w | X, y, \sigma^2) = \frac{p(y | X, w, \sigma^2) p(w)}{p(y | X)} = \frac{p(y, w | X)}{p(y | X)}$$

$$p(w | X, y, \sigma^2) = \frac{p(w, y | X, \sigma^2)}{p(y | X, \sigma^2)} \rightarrow \text{conditional probability}$$

$$= \frac{p(y | X, w, \sigma^2) p(w)}{p(y | X, \sigma^2)} \rightarrow \text{Bayes Theorem}$$

$$= \frac{p(y | X, w, \sigma^2) p(w)}{\int p(y | X, w, \sigma^2) p(w) dw} \rightarrow \text{marginal}$$

$$\int p(y|X, w, \sigma^2) p(w) dw = p(y|X, \sigma^2)$$

now, since we are integrating $p(w|y, X, \sigma^2)$, then $p(y|X, \sigma^2)$ will be a constant that normalize $p(w|y, X, \sigma^2)$

then, $p(w|y, X, \sigma^2)$ is likelihood prior
 $= p(y|X, w, \sigma^2) p(w)$

$$= N(X^T w, \sigma^2) N(\mu_0, \Sigma_0)$$

$$= \cancel{N(w, \sigma^2)}$$

$$= N(\mu_w, \Sigma_w)$$

$$\mu_w = \Sigma_w \left(\frac{1}{\sigma^2} X^T y + \Sigma_0^{-1} \mu_0 \right)$$

$$\Sigma_w = \left(\frac{1}{\sigma^2} X X^T + \Sigma_0^{-1} \right)^{-1}$$

Now, the posterior predictive PDF is:

~~$$\text{the marginalization of } f^* = p(f^* | w, X, y)$$~~

$$p(f^* | X, y) = \int p(f^* | w, X) p(w | X, y) dw \quad (\text{marginal})$$

$$p(f^* | w, X) = N(X^T w, \sigma^2) \rightarrow \text{by change of variable, get mean} = 0 \\ p(w | X, y) = \text{prior} = N(\mu_w, \Sigma_w)$$

↳ similar to sum of two Gaussians: $Z = X + Y$

$$\Rightarrow p(f^* | X, y) = N(X^T \mu_w, X^T \Sigma_w X)$$

$$\Rightarrow Z \sim N(\mu_w + \mu_y, \Sigma_w + \Sigma_y)$$

but b/c

$$F = X^T w \Rightarrow$$

$$F \sim N(X^T \mu_w, X^T \Sigma_w X)$$

Now, if we let $\phi(x)$ replace x as a basis function, to make the model more general, then we have

$$y = \phi(x)^T w + \varepsilon \rightarrow \tilde{y} = \bar{\phi}(x)^T w + \varepsilon, \quad x \in \mathbb{R}^D$$

$$\bar{\phi}(x) = [\phi(x_1), \phi(x_2), \dots, \phi(x_n)] \Leftarrow \begin{bmatrix} \phi(x_1) \\ \phi(x_2) \\ \vdots \\ \phi(x_n) \end{bmatrix} \quad \begin{array}{l} \text{each entry is a} \\ \text{function of vector} \\ x \text{ (input)} \end{array}$$

$$\in \mathbb{R}^{N \times n}$$

$$\Rightarrow \underbrace{w \in \mathbb{R}^{N \times 1}}_{\text{now, } N\text{-dimensional design matrix.}}$$

$$\hookrightarrow \begin{bmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{1n} \end{bmatrix}, x_i \in \mathbb{R}^D$$

Same results really ...

$$\text{let } p(w) = N(\mu_0, \Sigma_0) \quad (\text{prior})$$

~~$p(y | X, w, \sigma^2) = N(\phi(x)^T w, \sigma^2 I)$~~ $\quad (\text{likelihood})$

$$p(w | y, X, \sigma^2) = N(\mu_w, \Sigma_w) \quad (\text{posterior})$$

$$\mu_w = \Sigma_w \left(\frac{1}{\sigma^2} \bar{\phi}(x)^T y + \Sigma_0^{-1} \mu_0 \right)$$

$$\Sigma_w = \left(\frac{1}{\sigma^2} \bar{\phi}(x) \bar{\phi}(x)^T + \Sigma_0^{-1} \right)^{-1}$$

then

$$\Rightarrow p(y_{\text{new}} | X_{\text{new}}, X, y) = N\left(\bar{\phi}(x_{\text{new}})^T \mu_w, \bar{\phi}(x_{\text{new}})^T \Sigma_w \bar{\phi}(x_{\text{new}})\right)$$

If we let $K(X_1, X_2) = \bar{\phi}(X_1)^T \Sigma_w \bar{\phi}(X_2)$ then can re-write as:

$$p(y_{\text{new}} | X_{\text{new}}, X, y) = N\left(K(X_{\text{new}}, X)(K(X, X) + \sigma^2 I)^{-1} y, K(X_{\text{new}}, X_{\text{new}}) - K(X_{\text{new}}, X)(K(X, X) + \sigma^2 I)^{-1} K(X, X_{\text{new}})\right)$$

$$K(X_{\text{new}}, X_{\text{new}}) - K(X_{\text{new}}, X)(K(X, X) + \sigma^2 I)^{-1} K(X, X_{\text{new}})$$

Oh, so now we're actually looking at Gaussian processes
defined as a collection of random vars.

Let $f(x) \sim GP(m(x), K(x, x'))$ be the set of random variables
~~defining~~ representing a Gaussian Process.

notes: $E[f(x)] = m(x)$

$$\text{Var}[f(x)] = E[(f(x) - m(x))(f(x) - m(x))^\top] = K(x, x)$$

As a simple example,

$$f(x) = \phi(x)^\top w \quad w \sim N(0, \Sigma_p)$$

$$\Rightarrow E[f(x)] = E[\phi(x)^\top w] = \phi(x)^\top E[w] = 0$$

$$E[(f(x) - m(x))(f(x) - m(x))^\top]$$

$$= E[f(x) f(x)^\top] = E[\phi(x)^\top w w^\top \phi(x)] = \phi(x)^\top \Sigma_p \phi(x) = |K(x, x)|$$

If y is nugly training data, then $y = f(x) + \epsilon$, $\epsilon \sim N(0, \sigma_n^2 I)$

$$\Rightarrow y \sim N(0, K(x, x) + \sigma_n^2 I)$$

and for a joint distribution w/ new test points x^* :

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

$$\rightarrow f^* | y, x, x_* \sim N(\bar{f}_*, \text{cov}(f_*)) \quad \text{where} \quad (\text{per Gaussian conditional derivation})$$

$$\bar{f}_* = K(x_*, x) [K(x, x) + \sigma_n^2 I]^{-1} y$$

$$\text{cov}(f_*) = K(x_*, x_*) - K(x_*, x) [K(x, x) + \sigma_n^2 I]^{-1} K(x, x_*) \rightarrow$$

if we want to know y_x , then add $\sigma_n^2 \mathbf{I}$ to covariance

$$\Rightarrow y_x | y, X, x_* \sim N(\bar{f}_x, \text{cov}(f_x) + \sigma_n^2 \mathbf{I})$$

$K(x, x')$ is the "kernel"

an example of a kernel is the Squared Exponential kernel.

$$k(x_i, x_j) = \alpha \exp\left(-\frac{1}{2\ell^2} \|x_i - x_j\|_2^2\right), \quad \begin{array}{l} \text{amplitude } \alpha \\ \text{length-scale } \ell \\ \text{sum of infinite } \underbrace{\text{exponentials...}}_{\text{i.e. distance}} \\ \text{between two points} \rightarrow \text{"stationary"} \\ \text{also "isotropic" } \rightarrow \cancel{\text{absolute distance}} \\ \text{between points} \\ \hookrightarrow \text{the norm of distance} \end{array}$$

linear:

$$k(x, x') = \alpha x^T x' \quad \alpha > 0$$

polynomial

$$k(x, x') = \alpha (1 + x^T x')^r, \quad \alpha, r > 0$$

Cosine

$$k(x, x') = \alpha^2 \cos\left(\frac{2\pi}{\lambda^2} \|x - x'\|_2^2\right), \quad \alpha, \lambda > 0$$

Kernels can be combined to capture different aspects:

product $k(x, x') = k_A(x, x') k_B(x, x')$

scale $k(x, x') = g(x) k(x, x') g(x')$, arbitrary $g(x)$

sum $k(x, x') = k_A(x, x') + k_B(x, x')$

for our case, we have learned that is like:

periodic \leftarrow should we have multiple? ...

polynomial \leftarrow squared?

nugget \leftarrow SE (small length scale)

Oh so... once we have the desired K , we can simply plug it

along to get the multivariate distribution for new points

$$x_* \text{ (which can include } x) \rightarrow y_* | y, x, x_* \sim N(\bar{f}_{x_*}, \text{cov}(f_{x_*}) + \sigma_n^2 I)$$

to get desired $k(x, x')$ follow book example of RW.

also a parameter to estimate!

\rightarrow go make last one found.

$$k_1(x, x') = \theta_1^2 \exp\left(-\frac{(x-x')^2}{2\theta_1^2}\right) \rightarrow \begin{cases} \text{long term smooth} \\ \text{rising trend} \end{cases}$$

$$k_2(x, x') = \theta_3^2 \exp\left(-\frac{(x-x')^2}{2\theta_4^2} - \frac{2\sin^2(\pi(x-x'))}{\theta_5^2}\right) \rightarrow \begin{cases} \text{seasonal} \\ \text{variations} \end{cases}$$

allow seasonal

periodic component.
period = 1 (year). normally would
be $\frac{2\sin^2(\frac{(x-x')}{2})}{1^2}$

trend to decay away from periodicity

$$\theta_8, \theta_7 > 0$$

$$k_3(x, x') = \theta_6^2 \left(1 + \frac{(x-x')^2}{2\theta_8\theta_7^2} \right)^{\theta_8} \rightarrow (\text{small}) \text{ medium-term irregularities}$$

amplitude ↑ typical length-scale } rational quadratic term

shape parameter

For diffuseness of length-scales

(but radial quad. term better)

NOTE: could also use SE term for medium-term irregularities

$$k_4(x, x') = \theta_9^2 \exp\left(-\frac{(x-x')^2}{2\theta_{10}^2}\right) \rightarrow \text{short-term noise effects}$$

$$k_s(x, x') = \theta_{11}^2 I_n \rightarrow \text{independent noise of measurement...}$$

$n = \#$ training data

SD find best θ to fit own sample data (train)

$$p(y|x, \theta) \rightarrow y \sim N(\underbrace{y|D}_{\text{O-mean prior}}, C(\theta) + \theta_{11}^2 I_n), \quad C(\theta) = k_1(x, x') + k_2(x, x') + k_3(x, x') + k_4(x, x')$$

$C(\theta) = k(x, x')$

\rightarrow means have to subtract mean(y) before do analysis (!!)

$$\begin{aligned} \log(p(y|x, \theta)) &= \log N(0, C(\theta) + \theta_{11}^2 I_n) \\ &= \log \left[(2\pi)^{-\frac{n}{2}} |C(\theta) + \theta_{11}^2 I_n|^{-\frac{1}{2}} \exp \left(-\frac{1}{2} (C(\theta) + \theta_{11}^2 I_n)^{-1} y \right) \right] \end{aligned}$$

$$= -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log |C(\theta) + \theta_{11}^2 I_n| - \frac{1}{2} y^T (C(\theta) + \theta_{11}^2 I_n)^{-1} y$$

$n = \#$ of data points

gradient \rightarrow

Now find ~~maxima~~ maxima, using partial derivatives of w.r.t. θ_j :

$$\frac{\partial}{\partial \theta_j} \log p(y|x, \theta) = \frac{1}{2} y^T K^{-1} \frac{\partial K}{\partial \theta_j} K^{-1} y - \frac{1}{2} \text{tr}\left(K^{-1} \frac{\partial K}{\partial \theta_j}\right)$$

$$= \frac{1}{2} \text{tr}\left((\alpha \alpha^T - K^{-1}) \frac{\partial K}{\partial \theta_j}\right), \quad \alpha = K^{-1} y$$

$$\alpha^T = K^{-1} y^T K^{-1}$$

$$\text{will simplify a lot...?} \quad K = K(x, x')$$

optimum? "conjugate gradient optimizer" ...

\rightarrow Fréchet differentiable, no boundary cond. for $\theta_7, \theta_8 > 0$

$$\theta_1, \theta_5, \theta_4, \theta_2 \leq 0$$

use vpa solve? to get when all gradients = 0... yeah, cl goes on...

Want $\frac{\partial}{\partial \theta_j} \log p(y|x, \theta) = 0$ @ min/max \rightarrow iteratively solve & then compare values of multiple solutions \rightarrow from different starting points (local maxima)

\rightarrow once we have this, plug back in to find \bar{x}_k , $w_k(\bar{x}_k)$

~~draw from multivariate gaussian~~ \rightarrow draw from multivariate gaussian & plot!

Next step: 1) write up $K(x, x')$ in python & get expressions for gradients

2) ~~use~~ MATLAB solve for $\nabla_{\theta} \log p(y|x, \theta) = 0$

3) ~~and plot~~

(~~using gradient descent~~)

4) ~~and plot~~ plug in θ solved & plot

$$\frac{\partial h}{\partial \theta_j} = \frac{1}{\sigma^2} \begin{pmatrix} k_{11}(\theta) & k_{12}(\theta) & \cdots & k_{1n}(\theta) \\ k_{21}(\theta) & \ddots & \ddots & \vdots \\ \vdots & & & \\ k_{n1}(\theta) & \cdots & & k_{nn}(\theta) \end{pmatrix}$$

yup!

$$if K(x, x') = \sigma^2 \exp\left(-\frac{(x-x')^T(x-x')}{2\lambda^2}\right)$$

$$\frac{\partial h}{\partial \sigma} = 2\sigma \exp(-\dots)$$

$$\underline{\frac{\partial h}{\partial \lambda} = \sigma^2 \exp(-\dots) \frac{(x-x')^T(x-x')}{\lambda^3}}, \text{ yup}$$

~~if~~ $\frac{\partial}{\partial \theta} \log p(y|X, \theta) = 0$ when each derivative = 0 ...
also at other locs... have to calculate
for values of θ though.

So... we could generate a bunch of samples then
find local maxes & compare.
↳ need good starting points & enough

Gradient Descent for non-linear system auto-choose α (or 5 as it
 $0.0001 \rightarrow 13$ maybe)
logspace, 100 steps

~~1. steps:~~

1) ~~posterior not have structure~~

2) compute w/ nominal θ (b plot?) \rightarrow enforce rounding
eigenvalues to $\frac{6}{6}$

3) setup optimization for θ

(get eigenvectors & eigenvals,
round eigenvals to
check norms match $e^{-10} \dots$ should
be good..)

gradient of kernel derivation

$$h_1(x, x') = \theta_1^2 \exp\left(-\frac{(x-x')^2}{2\theta_2^2}\right)$$

$$\frac{\partial h_1}{\partial \theta_1} = 2\theta_1 \exp(-\dots)$$

$$\begin{aligned}\frac{\partial h_1}{\partial \theta_2} &= \theta_1^2 \exp(-\dots) (-2) \left(-\frac{(x-x')^2}{2\theta_2^3} \right) \\ &= 2\theta_1^2 \exp\left(-\frac{(x-x')^2}{2\theta_2^2}\right) \left(\frac{1}{2\theta_2^3} - \frac{(x-x')^2}{2\theta_2^3} \right) \\ &= 2\theta_1^2 \exp\left(-\frac{(x-x')^2}{2\theta_2^2}\right) \left(\frac{(x-x')^2}{8\theta_2^3} \right) \quad \checkmark\end{aligned}$$

$$h_2(x, x') = \theta_3^2 \exp\left(-\frac{(x-x')^2}{2\theta_4^2} - \frac{2\sin^2(\pi(x-x'))}{\theta_5^2}\right)$$

$$\frac{\partial h_2}{\partial \theta_3} = 2\theta_3 \exp(-\dots) \quad \checkmark$$

$$\frac{\partial h_2}{\partial \theta_4} = \theta_3^2 \exp(-\dots) \left(\frac{(x-x')^2}{\theta_4^3} \right) \quad \checkmark$$

$$\frac{\partial h_2}{\partial \theta_5} = \theta_3^2 \exp(-\dots) \left(\frac{2\sin^2(\pi(x-x'))}{\theta_5^3} \right) (2) \quad \checkmark$$

$$h_3(x, x') = \theta_6^2 \left(1 + \frac{(x-x')^2}{2\theta_7^2 \theta_8} \right)^{-\theta_8}$$

$$\frac{\partial h_3}{\partial \theta_6} = 2\theta_6 \left(\dots \right) \quad \checkmark$$

$$\begin{aligned} \frac{\partial h_3}{\partial \theta_7} &= \theta_6^2 \left(1 + \frac{(x-x')^2}{2\theta_7^2 \theta_8} \right)^{-\theta_8-1} \left(\frac{(x-x')^2}{2\theta_7^3 \theta_8} (-2) \right) \\ &= \theta_6^2 \left(1 + \frac{(x-x')^2}{2\theta_7^2 \theta_8} \right)^{-\theta_8-1} \left(\frac{(x-x')^2}{\theta_7^3} \right) \quad \checkmark \end{aligned}$$

$$\begin{aligned} \frac{\partial h_3}{\partial \theta_8} &= \theta_6^2 \frac{\partial}{\partial \theta_8} \left(1 + \frac{(x-x')^2}{2\theta_7^2 \theta_8} \right)^{-\theta_8} \\ &= \theta_6^2 \frac{\partial}{\partial \theta_8} \exp \left(\log \left(1 + \frac{(x-x')^2}{2\theta_7^2 \theta_8} \right)^{-\theta_8} \right) \end{aligned}$$

$$\text{ln}(u) = -c \log(\dots), \quad \frac{\partial}{\partial u} e^u = e^u, \quad \frac{\partial}{\partial \theta_8} = \frac{\partial}{\partial u} \frac{\partial u}{\partial \theta_8} \Rightarrow \frac{\partial e^u}{\partial \theta_8} = \frac{\partial e^u}{\partial u} \frac{\partial u}{\partial \theta_8}$$

$$\frac{\partial u}{\partial \theta_8} = -\log \left(\dots \right) + \left(\frac{1}{1 + \frac{(x-x')^2}{2\theta_7^2 \theta_8}} \right) \left(\frac{(x-x')^2}{2\theta_7^2 \theta_8} \right)$$

$$\frac{\partial h_3}{\partial \theta_8} = \theta_6^2 \left[\exp \left(\log \left(1 + \frac{(x-x')^2}{2\theta_7^2 \theta_8} \right)^{-\theta_8} \right) \right] \left[-\log \left(1 + \frac{(x-x')^2}{2\theta_7^2 \theta_8} \right) + \left(\frac{(x-x')^2}{2\theta_7^2 \theta_8 + (x-x')^2} \right) \right]$$

$$= \theta_6^2 \left[1 + \frac{(x-x')^2}{2\theta_7^2 \theta_8} \right]^{\theta_8} \left[\frac{(x-x')^2}{2\theta_7^2 \theta_8 + (x-x')^2} - \log \left(1 + \frac{(x-x')^2}{2\theta_7^2 \theta_8} \right) \right] \quad \checkmark$$

$$h_4 = \theta_q^2 \exp\left(-\frac{(x-x')^2}{2\theta_{10}^2}\right)$$

$$\frac{\partial h_4}{\partial \theta_q} = 2\theta_q \exp(\dots)$$

$$\frac{\partial h_4}{\partial \theta_{10}} = \theta_q^2 \exp(\dots) \left(\frac{(x-x')^2}{\theta_{10}^3} \right)$$

$$h_5 = \theta_{11}^2 I_m$$

$$\frac{\partial h_5}{\partial \theta_{11}} = 2\theta_{11}$$

$\log p(\theta)$. if $p(\theta)$ is a multivariate gaussian distribution, where each have their own mean & cov?

or it's like $p(\theta_1) + p(\theta_2) + p(\theta_3) + \dots + p(\theta_n)$ (so that
" " ← integrates to 1)

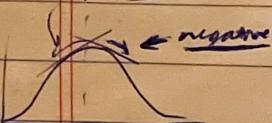
If Gaussian, then $p(\theta_i) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{1}{2} \frac{(\theta_i - \mu_i)^2}{\sigma_i^2}\right)$

$$\Rightarrow \log p(\theta_i) = -\frac{1}{2} \log(2\pi) + \log \sigma_i + \left(-\frac{1}{2} \frac{(\theta_i - \mu_i)^2}{\sigma_i^2}\right)$$

$$\frac{\partial \log p(\theta_i)}{\partial \theta_i} = 0 + 0 - \frac{1}{\sigma_i^2} \frac{(\theta_i - \mu_i)}{\sigma_i^2} (2)$$

$$= \frac{\mu_i - \theta_i}{\sigma_i^2}$$

stays positive



$$\text{if } b \sim (0, 1) \quad \mathbb{E}[bb^T] = I$$

$$b_2 \sim (0, 2) \quad , \quad \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \sim ?$$

$$\mathbb{E}\left[\begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right] = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\mathbb{E}\left[\begin{pmatrix} b_1 \\ b_2 \end{pmatrix}(b_1^T \ b_2)\right] = \mathbb{E}\begin{pmatrix} b_1 b_1^T & b_1 b_2^T \\ b_2 b_1^T & b_2 b_2^T \end{pmatrix} = \begin{pmatrix} \mathbb{E}[b_1 b_1^T] & \mathbb{E}[b_1 b_2^T] \\ \mathbb{E}[b_2 b_1^T] & \mathbb{E}[b_2 b_2^T] \end{pmatrix}$$

$$\mathbb{E}[b_1 b_2^T] = \int b_1 b_2^T f_{b_1 b_2}(b_1, b_2) db_1 db_2$$

Yeah... hard to tell/diffine w/o knowing distribution joint.
→ if features independent, then $= 0$

→ for extract basis function. $[1 \ x \ x^2 \ x^3 \ x^4]$

let them be independent w/ defined means + last get us in ball park

↳ normalize the dataset?? to make it easier for computing weights?