

## Técnicas para construir Kernels

Resolver por Lagrange para  $q=2$  y

$$\|M\|_2 = 1, M \geq 0$$

función  $\hat{B}(\sum_{j=1}^d \mu_j k(x_j, x_y))$

Restricciones:  $\|M\|_2^2 \leq 1$

Máximizan:

$$\mathcal{L}(M, \lambda) = \hat{B}(\sum_{j=1}^d \mu_j k(x_j, x_y)) + \lambda (1 - \sum_{j=1}^d \mu_j^2)$$

•  $\lambda$  = multiplicador lagrange asociado a la restricción  $q=2$ , es decir, norma  $l_2$ , implica  $\|M\|_2^2 = 1$

• Término  $\sum_{j=1}^d \mu_j^2$  representa norma  $l_2$  de los  $\mu_j$

Valores óptimos de  $\mu_j$ , lagrangiano respecto a  $\mu_j$  y  $\lambda$ .

$$\frac{\partial \mathcal{L}}{\partial \mu_j} = \frac{\partial \hat{B}}{\partial \mu_j} + \lambda (-2\mu_j) = 0$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = 1 - \sum_{j=1}^d \mu_j^2 = 0$$

Kernel.



## Hilbert Embeddings.

Si  $E_p \{ (K(x, x'))^{1/p} \} < \infty, \rightarrow M_p \in T.$

Operador  $T_p$

$$T_p f = E_p [K(x, x') f(x)]$$

$K(x, x')$  Kernel mide similitud  $x$  y  $x'$ .

$T_p$  operador acotado, función  $M_p$  tal  
que  $T_p f = C(f, M_p)$

Condición Inicial:

Si  $E_p [\sqrt{K(x, x')}] < \infty, M_p$  pertenece  
espacio funciones  $f$ .

$T_p f$

$$T_p f = E_p [K(x, x') f(x)]$$

$f(x)$  función actúa operador

Desigualdad Jensen cota superior:

$$|T_p f| = |E_p [K(x, x') f(x)]|$$

$$|T_p f| \leq E_p [K(x, x') |f(x)|]$$





$$\|T_P f\|_S = E_P [\|f(x)\| \sqrt{K(x, x')}]$$

$$T_P f = \langle f, M_P \rangle f.$$

$$M_P(x) = E_P [K(x, x')], \text{ } M_P. \text{ valor esperado Kernel,}$$

**GPR.** Gaussian Process Regression.

$$\text{Función: } f(x) \sim GP(m(x), K(x, x'))$$

donde  $m(x) = 0$  y  $K(x, x')$ , Kernel o función covarianza.

$$\{x, y\} \text{ con } y = f(x) + \epsilon, \text{ donde } \epsilon \sim N(0, \sigma^2).$$

es el ruido gaussiano, el objetivo es predecir el valor de  $f(x^*)$  para  $n$  nuevos puntos  $x^*$ .

La predicción  $f(x^*)$  dada la observación

y está dada por:

$$P(f^* | x, y, x^*) \sim N(M(x^*), \sigma^2(x^*))$$

donde.

$$M(x^*) = K(x^*, x) [K(x, x) + \sigma^2 I]^{-1} y$$

$$\sigma^2(x^*) = K(x^*, x^*) - K(x^*, x) [K(x, x) + \sigma^2 I]^{-1} K(x, x^*)$$

$$+ \sigma^2 I]^{-1} K(x, x^*).$$



Aquí,  $K(X, X)$  es la matriz de covarianza construida con los puntos  $X$ .

Optimización:

hiperparámetros, longitud escala  $L$  y variancia  $\sigma^2$  maximizando logaritmo de la verosimilitud marginal:

$$\log p(y|X) = -\frac{1}{2} y^T [K(X, X) + \sigma^2 I]^{-1} y - \frac{1}{2} \log |K(X, X) + \sigma^2 I| - \frac{n}{2} \log 2\pi$$

Optimiza respecto a los hiperparámetros del kernel usando métodos como L-BFGS o Adam.

**GPC** Gaussian Process Classification.

función  $p(y=1|f(x)) = \phi(f(x))$ .

$\phi$  función sigmoide o probit.

Optimización:

Métodos de aproximación.

Laplace. Aproxima verosimilitud marginal alrededor del valor máximo posterior.

Expectation Propagation (EP): Aproximación iterativa para computar la verosimilitud marginal.



## VGP. Variational Gaussian Process

función.: Aproximación variacional para conjunto de datos grande.

Aproximar  $p(f | X, y)$  con distribución variacional  $q(f)$ .

Evidencia variacional inferior. Cota inferior a la log-verosimilitud marginal:

$$\text{ELBO} = \mathbb{E}_q[\log p(y|f)] - \text{KL}(q(f) \| p(f))$$

KL es divergencia kullback-leibler entre aproximación variacional  $q(f)$  y la verdadera distribución posterior  $p(f)$ .

Optimización:

Elbo se maximiza respecto parámetros variacionales y los hiperparámetros del kernel. El gradiente de la Elbo se calcula mediante técnicas como la diferenciación automática y se optimiza con métodos de gradiente.

## SGPR. Sparse Gaussian Process Regression.

función.

Conjunto puntos ~~inducción~~ inducción  $Z$ , que se aproximan a posterior.



Reducir complejidad computacional de  $O(n^3)$  a  $O(m^2n)$ , donde  $m \ll n$ .

ELBO en SVPR es:

$$\text{ELBO} = \mathbb{E}_q[\log p(y|f)] - \text{KL}(q(f) \| p(f|z))$$

Optimización:

Maximizar puntos de inducción  $z$  y los hiperparámetros del Kernel maximizando ELBO.

Se realiza con optimización de gradiente, (Adam).

**SVGP.** Sparse Variational Gaussian Process  
Generaliza a clasificación y Regresión,  
manteniendo estructura y puntos de  
Inducción para hacer aproximación varia-  
cional más eficiente.

$$\text{ELBO} = \mathbb{E}_q[\log p(y|f)] - \text{KL}(q(u) \| p(u)).$$

$q(u)$  es distribución variacional sobre puntos de inducción.

Optimización: ELBO respecto puntos inducción e hiperparámetros del Kernel.