

```
In [ ]: import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import load_breast_cancer
import torch
import torch.optim as optim
from torch.optim import Adam
import torch.nn.functional as F
from torch.autograd.functional import hessian
from torch.distributions.multivariate_normal import MultivariateNormal
import seaborn as sns
import io
import base64
```

Instruções gerais: Sua submissão deve conter:

1. Um "ipynb" com seu código e as soluções dos problemas
2. Uma versão pdf do ipynb

Caso você opte por resolver as questões de "papel e caneta" em um editor de $LAT_{E}X$ externo, o inclua no final da versão pdf do 'ipynb'--- submetendo um único pdf.

Trabalho de casa 05: Processos Gaussianos para regressão

1. Durante a aula, discutimos como construir uma priori GP e o formato da posteriori preditiva para problemas de regressão com verossimilhança Gaussiana (com média definida pelo GP). O código abaixo cria um GP com kernel exponencial quadrático, mostra a priori preditiva e a posteriori preditiva. Experimente com o código e comente a influência de ambos os parâmetros do kernel exponencial quadrático, tanto na priori preditiva quanto na posteriori preditiva. Nos gráficos gerados, os pontos vermelhos são observações, as curvas sólidas azuis são as médias das preditivas e o sombreado denota +- um desvio padrão.

```
In [ ]: SEED = 42
np.random.seed(SEED)
s2 = 1e-04 # variância observacional

# Kernel Exponencial Quadrático
def rbf_kernel(x1, x2, gamma=10.0, c=1.0):
    assert(gamma>0)
    assert(c>0)
    return (-gamma*(torch.cdist(x1, x2)**2)).exp()*c

# Posterior Predictive
def posterior_pred(x, xt, yt, gamma=10.0, c=1.0, sq=s2):
    Kxxt = rbf_kernel(x, xt, gamma, c)
    Kxt = rbf_kernel(xt, xt, gamma, c) + torch.eye(xt.shape[0])*sq
    Kinv = torch.linalg.inv(Kxt)
    Kxx = rbf_kernel(x, x, gamma, c)

    mu = Kxxt @ Kinv @ yt
    cov = Kxx - Kxxt @ Kinv @ Kxt.T
    return mu, cov

# Plot Gaussian Process
def plot_GP(gamma: list, c: list) -> None:
    n_params = len(gamma)
    fig, axes = plt.subplots(int(n_params/2), 4, figsize=(20, 2 * n_params))
    axes = axes.reshape(-1, 2)
    # Loop para cada parâmetro gamma e c
    for gamma, c, ax in zip(gamma, c, axes):
        x = torch.linspace(-1, 1, 100).reshape(-1, 1) # Criação de um vetor x com 100 pontos entre -1 e 1
        # Cálculo da matriz de covariância K
        K = rbf_kernel(x, x, gamma, c) + torch.eye(x.shape[0])*s2
        # Criação de um vetor de médias mu com zeros
        mu = torch.zeros_like(x)
        # Definição dos pontos de treino
        xtrain = torch.tensor([-0.5, 0.0, 0.75])[:, None]
        ytrain = torch.tensor([-1.5, 1.0, 0.5])[:, None]

        # Plot da média e intervalo de confiança para o prior
        ax[0].plot(x, mu)
        ax[0].set_xticks(xtrain.flatten())
        ax[0].fill_between(x.flatten(), mu.flatten()-K.diag(), mu.flatten()+K.diag(), alpha=0.5)
        ax[0].set_xlim([-1, 1])
        ax[0].set_ylim([-5, 5])
        ax[0].set_title(f'GP prior (gamma: {gamma}, c: {c})')

        # Cálculo da média e covariância posterior
        post_mu, post_cov = posterior_pred(x, xtrain, ytrain, gamma, c)
        # Plot da média e intervalo de confiança para o posterior
        ax[1].plot(x, post_mu)
        ax[1].set_xticks(xtrain.flatten())
        ax[1].fill_between(x.flatten(), post_mu.flatten()-post_cov.diag(), post_mu.flatten()+post_cov.diag(), alpha=0.5)
        ax[1].scatter(xtrain, ytrain, color='red', zorder=5)

        ax[1].set_xlim([-1, 1])
        ax[1].set_ylim([-5, 5])
```

```

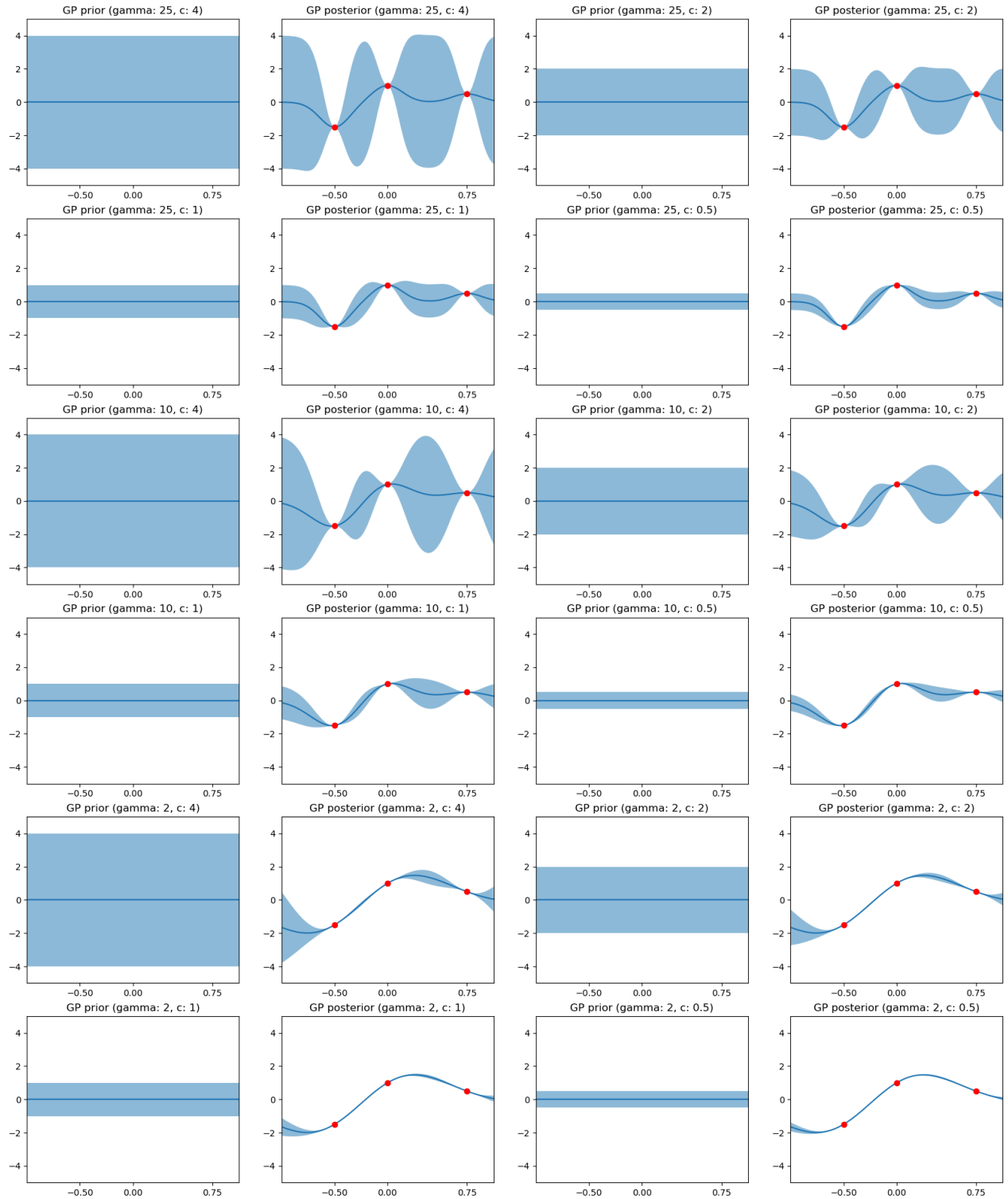
    axs[1].set_title(f'GP posterior (gamma: {gamma}, c: {c})')
    plt.show()

```

```

In [ ]: plot_GP([25, 25, 25, 25, 10, 10, 10, 10, 2, 2, 2, 2], [4, 2, 1, 0.5, 4, 2, 1, 0.5, 4, 2, 1, 0.5])

```



Defino $\gamma = 25, 10, 2$ e $c = 4, 2, 1, 0.5$. Note que:

Na GP Prior, um γ menor resulta em uma função média suave e previsível, enquanto um c menor restringe a variabilidade vertical das funções.

Na GP Posterior, a incerteza do modelo é reduzida perto dos pontos de dados observados, com o efeito de um γ pequeno ainda evidente na suavidade da função média.

- O pequeno valor de γ resultou em uma função suave, tanto na priori quanto na posteriori.
- O pequeno valor de c restringiu a variabilidade das funções sampleadas do GP, mantendo as previsões mais próximas à média e resultando em uma incerteza reduzida (sombreado estreito).

2. Durante a aula, discutimos como escolher os hiper-parâmetros do nosso GP. Estime os parâmetros ótimos para os dados carregados abaixo (acredite, é isso que o código faz). Reporte a evidência obtida e faça um plot similar ao acima. Para o dado de teste, reporte a i) log verossimilhança e ii) o MSE com relação à média. Em caso de dúvidas, recorra a nota de aula e o link adicionado no eclass.

```
data = np.load(io.BytesIO(ase64.b85decode('P)h@6awAK2mk;8Ap0$!k1tLx00Sv000XB6aaK' VQfQc(STlGcGw?r0HGzCw000000D)Vn00000G-JFJ_});!0~YzY=9g==Y&XhJ79<?P+Lz&d|0HoP>&SM')2NI XZ(Ac!XAF8fg9xVzSb1PlJNq0ihTRBQA_Q_*23-qcCzMESCGRBV!9&(A)hLg<_iy14>x0ld$%0%-J21(ae!R+&zB5xiP_e970F;_<;vgd|KI T18'SS'&@8X)%92ftnhkr)}!{KI3SqaztHTO%1?F#czrlL<L1=>jXSP9_-ra#1EFBN&DUp#j93U}z=UqW<3m;jT?g?mc%DzX*@x|yc(1GbZBQC_u|rk04V2BBL^Y+2'2>OV8sd>fT|N5_c9=j0c<;Ux)khM$ssKK0%)5=2f&UT^G<B|Ttrn^Vx M7!hh- (FZfcVtH+vFY0ZHVTFU7#>wvZ_WXhYwJL>hZevBPdqT^trs'DhO5DX6wZ>'!v-Va->Dtd'Ede=kXs3rEhnf py+{YH_yeuY Z11;a)XkTs17Za04Q26FSDQ3_OfxLMq@b(L+S2Ghtyia'8XS&mna<{tu'x8IS3?1fudzleFnF10Qq'+kbUaeQpyT1EdkEdic>IhkEac-Msa!K!)ShvPHNAZV)_Rs_ju9td17?gYFSad-g_Y_arJb^!hCNgeWsl1jg6fk9PCRPWPpLPu>+<k;l1>Jd)FHHX1(ZWGZ'C;jap_g;LhbT)#a(VMt#9tJ9t)!vXuq6qn={f|sqnfICf9an-Xxcex}zY}Cs9U0wi6SLmp$Vj62^s)Hnt9+sTc3wi*r;gfKt%AR%64EI#Z8d7FHaQU@N'frb1=1vY_4L7yh-uq#wd53_qdQs<1TeGe<iat;hn&RT6P7JRVAS10qWb1lvnm?GPVU0?7Ng9cDLKWbVqC&FAmxT1kAuZ^F65;Z2dG67)Esds=y8P^UoGS^WFK'wcvb>Yq<XC09u4=2f&bW8_v0_l&n<1 q oq=k';7=Hp|jX10uDMqq_k'Dmgte@AwPyVsqt09Z4Fo1>1>I>ja9VXWtYJCLb<_50HrcCzBs71i>2co'tJoQ2aQa'N(u'b=2&rY#1ffUM&Q*1;RC6!5@Yw1Yxmo7=qAT7Y4vhOHYq#!de*%FTX3j}%0ZX}3~rwKGbt4aw^j#Fnz!28|wixCgxmhJGG<AU+R5ARn<n6GTOWiXp(Ia'P)h>@6awAK2mk;8Apk!*Q7DH400Sv000XB6aaK' VQfQc(c'j~ncw>wr0HGzCw000000EYt000000G^KSGaCjRFKy_8T*7DENx=P);51LtqPkws_[k~h(iOYe'O?L?cl?BCWTWE8&I-f)a+>|5NB+ODErWN$XRgdiehYivqv1IVv~>?XX@;1;3Ep&@QL%_<)K)NK9KC==CgjnIDQ?j?HH97zlNoLPwdCl!fatuCYLyrDrTo_';'k!aoag+g5Fl$Emw-u'g GXB-_YPypHcmCWQR1ZtoCB(vOzFCYKAgJDNm@54duk3PdbIm0DDZGC&UHwh'|{pJ'_DdgZi>?DzubqG7CqCpgpGEf5sg&cODrk&S'h&cASfrB($S$FR@kY=?lb<(Gib8=If@uL-N#|GT)&W>R30q<aG>0^'8IsfaOu5M&p|-|t>h-hh_-+8izM_D3Rsc<o7_Dct(U@d$C!3=QvkbpQrPTP^GjTUhg)d'GN9Tr5<P>VBG&gCUW7xhvxb3mHBNNC}<uWJZ@Y)r6?}N5p(GF>s)n|OT9Tr9DDBr1Pnp>8ITSfLFDcv*_d9p3ZghY&Z-mG%-Mjouw=M&Rjk&a<3?38&oeozf'&[rm1jZQ4AH&lgSPUEW40?>BadCI?#4qOKRbY_l=@0mmIEp2tpQ>3ChZDyb91-wd4mrAX?>+AI=->Xo|y;$&(-NTKSFYDe#sIkM-1#3Z1?sn&MRuD30JAcoox'l8q88>=k7BU19H67Y=0-4awMYifd-Cg70_CAWAcc5019;1I9a=(4diTv{w;hb51Yhc3eIID)qg{lq<otU+L0rw}WgcdXu?pzZ6g3h{1LfK-8uf-KWK9;7X#=#2MTqTCfmHBoT$aauw04|zbj;h1oxmb@0}G:841YXz;dw9Y_QQJ1}'&GX%lt0Qt5SoL+X*2^|m4wUDAs^C#(img;=r7PSeYu70^DLe$uj_xy&)-u-j@vwXvj'X)G^x&s'rMAP?#!2_AzXuc->+?SAJj1Di8(c{&awpWPX@*&Gf&IA*5@7RI9X&UPH?G6Y-eKB5821<Y?yVPK0$S5B%3Y--_e4JfQ&A<-REUd?>?Rwnb;9a@&'whnxLtm(DAOA0EqH&lgSPUEW40?>F-FotH2CZVvGoC2' jxrXxz2bQ' #a@?r7#mq6{*vaGLMW4+XkIP5SVG!1@XVOEEq5dj7%e4JfQ&d8tkC'Rdy'YUNBYrn~*V'RM*14P)h>@6awAK2mk;8Apn6>1aKG%001Bm000UA6aaK(b97%=E^csnORRVHAP@im00007zzLY00000omcx~j%6R8w_g=49Or^hm)ny%_{j!$)TxvfSIKK'skiZWbwC_J_tk_fVx~deEEDY)Xuy5mfX<_0+4+CqGefGddtdkg(TpmchwAS_j0;-xt^vmOS;u&'H'dC1;e6%;{QEL_sYV)1dBqe21@?Q(z#&H!Ve#8c')VIBW_XEDto^pp^KYhn~p<?>KH0?>4409_TmQ@Be(Y?HQh?>@C?>X?>C5J5_-6Hu)-BetN%=(huzc^R5D^X9dyjr=Q0ZtI|L(8<S'H'wbKQ8H*1TxZG3S9KU)cq@VYGkeOBLL'5e3N&G(64)ZugVI7j9tvj<qZK_HN)'>#7gXyBUD(L>jrujKcKBkdQw)Ydx3Z(2QrZranT'=w2bz8^OFHEgxsdd))&($+4=$_ $1J?>oj=u6k)hx7?>zeXx(R&yJ4q%)#zdc9~>0Tt~CyVD=8p_Ed8sgsvrs@?TAq14;br=D^> <jyu3nzp-13T^-dIX_mQ$E!$HG&fZnai-18n{<x<=>'hbw|d9-IDiv-;KXwVda*JOYNDEyvXbk_KT87w#@UTSF_WIFPhXKrO@|Tt~C6TFPDdsWndPLU7^@ZJPBJL7e9S)ZA=?or^YERMcMGx_xrXNH'<}k00Yi+xk'(UANRQIO);?J0Z-gsT30' ZVIZgnst'>@C(g+08#~YFPYHrC>5ydr' $6W151r&Y'nUWh2lOt9wa(wtLn'=chPiOh==SQdb16%'vSHZ1jje?>Q&*Q'HRsI8dFY1&)f_v|_G_Iovw^tiY>67By)Ok'hJ?>qCTdxw(gcw*NrtmdlyZ-yvz^J>Nx^?>vGBMHOGLCk>?y6?>OTII1$dIdEv=L9^%QR<YDNMcnu^ZtLK(C|M#I)'wl$g$9jiF19^9Z-8TFXM|CVhq03{Q@%|PfB;#09&SdrFpc4Zezi~~~?Nty(yzjbY&A;HXZ+f_Zu90rTc%(DxDotUHC6S^;uQ_S_ *T2?}(=XWNY6m7Y81;< <V'd-aB<3C1<1j% yxpSBK~'>^*6di+g6WV|1X3+t>ZrFE)I87M3%1Ztc0z3s31~*r-QH+JdnU0^v|yqCB!_vuItCZH1rHeM-KXw5P'ySxM'N5#5j4JCidev|_10J4PAPam0|nr#JFfbviBg@>SXhj>ob>Y>OuO(&{j(LNG^z0BQ_6_43x<Z~7~itAdmN4Q?>0+k0;|zt>b>7roX)pSD$eETs^je_c~?J$pYX&3Er^vMaJ'@6WGYBW8Ini#01WRXSD>+T+>|KGI8H?>[8YIbg|_%3C)ERJj^rTottj;LrLN8^#keECE(Lq^>ZVzhB2+>fo7MUZC2Qd4r?>?MCR)(3ft?PP4x@K>PA>r>}>K'Vrn^tQ1?>H1cL(6G@6A8JicV_r'ePYpx&Za2VU>ZOPZ|_0_$#Px<D<#7fwCPmlQ0?>7HopJX;|U;Gzy?>K)bRezo7gouC4D3#ENqmN38L<EvndvG;GAMI@mM^<K'>q4T(pCS-E~>1050zxu&B;J0IT^Ymv_+on78z|;B0TJZGewjt~tiho@TLIGUN4UAWy93Hx=>EY0tDmrke5v6><=FnxM';2>eYZz~Z-3w2g811m^>B'Zj'>'PFQBQ^+4B|a^yI)e0*k(8NdQWTkw>JH^WQ|nCuv^3T|<LdkbBCiviue3cyP*8(^5|D1D)NtMmwUaiV2gspjz~5IJHi$MJTnh1vzX|tr_|Vk^udbCSBV(g'3Pl2Lq4!dn'FuokJoqpTDZ++fM40)|JE^E&GSWBSTX^!mc^4QVPefLkL+NQ01Mnf-gxg20Q<^<IQ^AQI62ZfIDQY~6T(Usx=f<Q|DiZFzRp9m#Dpz2t+iR&_1G|>Q0(QZ#GA&I)F10dpmtFX@+=+X7U0|_mn-XZ{<EEno8NLX5>Kp5=>=pqEtEA3<S-NoqxX7YBI_5M=I5E&A4|3nT_LXLS;qmrJM;GK^>AMNH+E2z^>K&E&T~FLJ'$YNYjmVmV=VgrZtemT~>+)%GIG=J52hMg>yyP;1D~L69AbpwCfi0dg2|xUV)4u1Te;=QRu(k3F!49W~>~?x7G$acXnU6$^NPBFWE$Ak^>nZrTO>Op8PmZ-J'>';Z(R=)Aq=I_i_4*s-2IHIZ0_<45i^D=hLYSHUjbejhMz(qzbKU1*cIN&';Gubk|AKkf70+>%TX8+z+hbcB4_Tps58' Rty%>+UnMu4p=NQ=dAz>T06do0'>?)SbhI08z{U2Wx9ptYpUt+nvVS_fwQ0Yua967*d7X&??>?Tnb65vqM-wa$TNiV2uqPGsxcsfn~2TXJk>I6U44d68NMNKjpornaGFkn&4ofKc(drXy)HFfv'AJI&GMUE84ec2CMWH3UtlP(I=0RMR?>R%>ZdQPar8<_wcgcKt3ZV!7L7LG65)7&CLBbxpf?YYK6%|jvh>vYVV3crncyGB8%|VUa(BV-6y2ou1wIw=diE1fo'>'i8SF9Lxfwnng!'<^*GLR3Q6Su-qj&&six;p}Jb@d;8weARK(urO)#4fs7VjWIjCNvdCHRV'5Fy?|ns^LXwEzZ=-de~2bXvdfOBPIYj
```

```
'0000X06#iWD2D?80H6Z^01E&B00000000000Du9&0{{SYa$#w1UwJNWaCuNm0Rj{Q6aWAK2mk;8Apn6>1aKG%001Bm000UA00000000004jiga-fsby*jN'
'Usx_~aCuNm0Rj{Q6aWAK2mk;8ApnKqayn=a001Bm000UA000000000004ji*bx8#bY*jNUwJNWaCuNm1qJ{B000C410Vay004X;00000'
)))

train_X, train_y = data['train_X'], data['train_y']
test_X, test_y = data['test_X'], data['test_y']

# Converte os dados para tensores
X_train = torch.tensor(train_X, dtype=torch.float32)
y_train = torch.tensor(train_y, dtype=torch.float32)
X_test = torch.tensor(test_X, dtype=torch.float32)
y_test = torch.tensor(test_y, dtype=torch.float32)
```

```
In [ ]: # Função para calcular a Log-verossimilhança marginal negativa
def GP_log_likelihood(x, y, gamma, c, sigma_noise_sq):
    K = rbf_kernel(x, x, gamma, c) + sigma_noise_sq * torch.eye(x.size(0))
    L = torch.cholesky(K)
    return -0.5 * (y.t() @ torch.cholesky_solve(y, L) + torch.logdet(L) + x.size(0) * torch.log(torch.tensor(2 * torch.pi)))

# Função para calcular o erro quadrático médio
def mse(y_pred, y):
    return torch.mean((y_pred - y) ** 2)
```

```
In [ ]: # Inicialização de hiperparâmetros
gamma = torch.tensor([10.0], requires_grad=True)
c = torch.tensor([1.0], requires_grad=True)
sigma_noise_sq = torch.tensor([0.1], requires_grad=True) # Ruído observacional
num_iterations = 1000 # Número de iterações
```

```
In [ ]: # Inicializa o otimizador
optimizer = Adam([gamma, c, sigma_noise_sq], lr=0.01)

# Loop de otimização
for i in range(num_iterations): # num_iterations é o número de etapas de otimização que você deseja executar
    optimizer.zero_grad()
    log_likelihood = -GP_log_likelihood(X_train, y_train, gamma, c, sigma_noise_sq) # Função para calcular a Log-verossimilhança
    log_likelihood.backward()
    optimizer.step()

    if (i + 1) % 100 == 0:
        print(f"Iteração {i+1}: Log Likelihood = {log_likelihood.item()}")

# Após a otimização, use os valores otimizados de gamma e c para calcular a média e a covariância do GP posterior e fazer prev
x = torch.linspace(-1, 1, 100)[: , None]
mu = torch.zeros_like(x)
K = rbf_kernel(x, x, gamma, c) + torch.eye(x.shape[0])*sigma_noise_sq
```

C:\Users\iaram\AppData\Local\Temp\ipykernel_21624\565260635.py:4: UserWarning: torch.cholesky is deprecated in favor of torch.linalg.cholesky and will be removed in a future PyTorch release.

```
L = torch.cholesky(A)
should be replaced with
L = torch.linalg.cholesky(A)
and
U = torch.cholesky(A, upper=True)
should be replaced with
U = torch.linalg.cholesky(A).mH
This transform will produce equivalent results for all valid (symmetric positive definite) inputs. (Triggered internally at
..\aten\src\ATen\native\BatchLinearAlgebra.cpp:1703.)
L = torch.cholesky(K)
```

```
Iteração 100: Log Likelihood = 35.39397430419922
Iteração 200: Log Likelihood = 35.08833312988281
Iteração 300: Log Likelihood = 34.84060287475586
Iteração 400: Log Likelihood = 34.6528205871582
Iteração 500: Log Likelihood = 34.51133728027344
Iteração 600: Log Likelihood = 34.40361022949219
Iteração 700: Log Likelihood = 34.32072448730469
Iteração 800: Log Likelihood = 34.256385803222656
Iteração 900: Log Likelihood = 34.20603561401367
Iteração 1000: Log Likelihood = 34.166404724121094
```

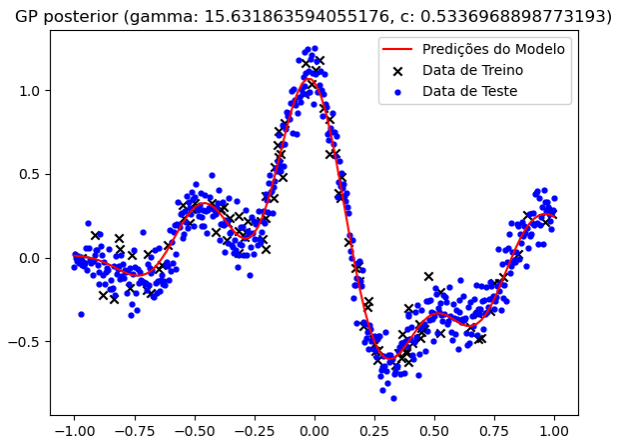
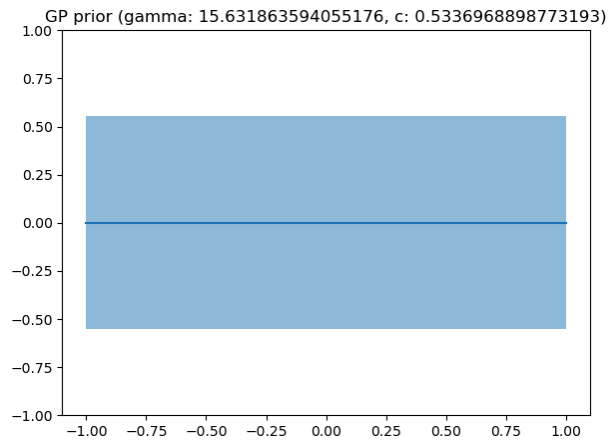
```
In [ ]: # Plotting
fig, axs = plt.subplots(1, 2, figsize=(15, 5)) # Create two subplots
# Parâmetros otimizados
gamma_opt, c_opt, sigma_noise_sq_opt = gamma.item(), c.item(), sigma_noise_sq.item()

# GP Prior
axs[0].plot(x.detach().numpy(), mu.detach().numpy())
axs[0].fill_between(x.flatten().detach().numpy(), mu.flatten().detach().numpy()-K.diag().detach().numpy(), mu.flatten().detach().numpy()+K.diag().detach().numpy(), alpha=0.7)
axs[0].set_ylim([-1, 1])
axs[0].set_title(f'GP prior (gamma: {gamma_opt}, c: {c_opt})')

post_mu, post_cov = posterior_pred(x, X_train, y_train, gamma=gamma_opt, c=c_opt, sq=sigma_noise_sq_opt)

# GP Posterior
axs[1].plot(x, post_mu, 'r', label='Predições do Modelo')
axs[1].fill_between(x.flatten(), post_mu.flatten()-post_cov.diag(), post_mu.flatten()+post_cov.diag(), alpha=0.7)
axs[1].scatter(X_train, y_train, color='black', label='Data de Treino', marker='x')
axs[1].scatter(X_test, y_test, color='blue', label='Data de Teste', s=12)
axs[1].set_title(f'GP posterior (gamma: {gamma_opt}, c: {c_opt})')
```

```
axs[1].legend()
plt.show()
```



```
In [ ]: x = torch.linspace(-1, 1, 500)[: , None]
# Calcula a Log-verossimilhança e o MSE nos dados de teste
test_log_likelihood = -GP_log_likelihood(X_test, y_test, gamma_opt, c_opt, sigma_noise_sq_opt)
post_mu, post_cov = posterior_pred(x, X_test, y_test, gamma=gamma_opt, c=c_opt, sq=sigma_noise_sq_opt)
```

```
In [ ]: mse_ = mse(post_mu, y_test)
print(f"MSE: {mse_}")
print(f"Test Log Likelihood: {test_log_likelihood.item()}")
```

MSE: 0.00912289135158062
Test Log Likelihood: 109.20968627929688

Exercício de "papel e caneta"

1. Na nota de aula, derivamos a posteriori preditiva $p(y_*|x_*, x_1, y_1, \dots, x_N, y_N)$. Por simplicidade, deduzimos a priori preditiva $p(y_*, y_1, \dots, y_N|x_*, x_1, \dots, x_N)$ e as condicionamos nas saídas y_1, \dots, y_N observadas no conjunto de treino. No entanto, também é possível obter o mesmo resultado calculando a posteriori $p(f_*, f_1, \dots, f_N|x_*, x_1, y_1, \dots, x_N, y_N)$ e, então, calculando o valor esperado de $p(y_*|x_*, f_*)$ sob essa posteriori. Deduza novamente a posteriori preditiva seguindo esse outro procedimento.

(Dica: você pode calcular a conjunta $p(f^*, f_1, \dots, f_N, y_1, \dots, y_N|x_*, x_1, \dots, x_N)$, que também será Gaussiana.)

Resposta:

Para calcular a posteriori preditiva, vamos considerar a distribuição conjunta de $f^*, f_1, \dots, f_N, y_1, \dots, y_N$, tal que a média é 0, pois assumimos que o GP é centrado. Sabendo que as observações y são relacionadas com f através de uma relação de ruído, normalmente assumida como normalmente distribuída, temos:

$$y_* = f_* + \epsilon$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2)$$

Para encontrar a matriz de covariância, iremos calcular a covariância entre as funções latentes e as saídas para os 3 possíveis casos. Seja $i, j \in \{1, \dots, N\}$, temos:

$$\text{Cov}(f_i, f_j) = k(x_i, x_j) \quad (1)$$

$$\text{Cov}(f_i, y_j) = \text{Cov}(f_i, f_j + \epsilon_j) = \text{Cov}(f_i, f_j) + \underbrace{\text{Cov}(f_i, \epsilon_j)}_0 = k(x_i, x_j) \quad (2)$$

$$\text{Cov}(y_i, y_j) = \text{Cov}(f_i + \epsilon_i, f_j + \epsilon_j) = \text{Cov}(f_i, f_j) + \underbrace{\text{Cov}(f_i, \epsilon_j)}_0 + \underbrace{\text{Cov}(\epsilon_i, f_j)}_0 + \underbrace{\text{Cov}(\epsilon_i, \epsilon_j)}_{(*)} \quad (3)$$

$$= k(x_i, x_j) + (*) \quad (4)$$

$$(*) = \begin{cases} \sigma^2, & \text{se } i = j \\ 0, & \text{se } i \neq j \end{cases} \quad (5)$$

Portanto, a distribuição conjunta é dada por:

$$\begin{bmatrix} f^* \\ f_1 \\ \vdots \\ f_N \\ y_1 \\ \vdots \\ y_N \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{k}_{**} & \mathbf{k}_{*X} & \mathbf{k}_{*X} \\ \mathbf{k}_{X*} & \mathbf{k}_{XX} & \mathbf{k}_{XX} \\ \mathbf{k}_{X*} & \mathbf{k}_{XX} & \mathbf{k}_{XX} + \sigma^2 \mathbf{I} \end{bmatrix} \right)$$

Onde:

- \mathbf{k}_{**} é a variância do novo ponto x_*
- \mathbf{k}_{*X} é o vetor de covariâncias entre o novo ponto x_* e os pontos de treinamento $X: x_1, \dots, x_N$
- \mathbf{k}_{X*} é a transposta de \mathbf{k}_{*X}
- \mathbf{k}_{XX} é o vetor de covariância entre os pontos de treinamento.

Vamos agrupar a matriz de covariância em 4 blocos da seguinte forma:

$$\Sigma_{1,1} = \begin{bmatrix} \mathbf{k}_{**} & \mathbf{k}_{*X} \\ \mathbf{k}_{X*} & \mathbf{k}_{XX} \end{bmatrix} \quad (6)$$

$$\Sigma_{1,2} = \begin{bmatrix} \mathbf{k}_{*X} \\ \mathbf{k}_{XX} \end{bmatrix} \quad (7)$$

$$\Sigma_{2,1} = [\mathbf{k}_{X*} \quad \mathbf{k}_{XX}] \quad (8)$$

$$\Sigma_{2,2} = [\mathbf{k}_{XX} \quad \mathbf{k}_{XX} + \sigma^2 \mathbf{I}] \quad (9)$$

Utilizaremos a propriedade condicional das distribuições Gaussianas para calcular a posteriori de \mathbf{f} condicional às observações \mathbf{y} :

$$\mu_{\mathbf{f}|\mathbf{y}} = \mu_1 + \Sigma_{1,2} \Sigma_{2,2}^{-1} (\mathbf{y} - \mu_2) \quad (10)$$

$$= \Sigma_{1,2} \Sigma_{2,2}^{-1} \mathbf{y} \quad (11)$$

$$\Sigma_{\mathbf{f}|\mathbf{y}} = \Sigma_{1,1} - \Sigma_{1,2} \Sigma_{2,2}^{-1} \Sigma_{2,1}$$

Ou seja,

$$\begin{bmatrix} f^* \\ f_1 \\ \vdots \\ f_N \end{bmatrix} \mid \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \sim \mathcal{N}(\mu_{\mathbf{f}|\mathbf{y}}, \Sigma_{\mathbf{f}|\mathbf{y}})$$

Usando a propriedade de gaussiana, temos que a marginal de f^* é:

$$f^* | y_1, \dots, y_N \sim \mathcal{N}(\mathbf{k}_{*X} (\mathbf{k}_{XX} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}, \quad \mathbf{k}_{**} - \mathbf{k}_{*X} (\mathbf{k}_{XX} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_{X*})$$

Finalmente, a predição de $y^* = f^* + \epsilon$ onde $\epsilon \sim \mathcal{N}(0, \sigma^2)$, resulta em:

$$y^* | x^*, y_1, \dots, y_N \sim \mathcal{N}(\mathbf{k}_{*X} (\mathbf{k}_{XX} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}, \quad \mathbf{k}_{**} - \mathbf{k}_{*X} (\mathbf{k}_{XX} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_{X*} + \sigma^2)$$

2. Quando trocamos a verossimilhança Gaussiana por uma Bernoulli (i.e., no caso de classificação binária), a posteriori para nosso GP não possui fórmula fechada. Mais especificamente, a verossimilhança para esse modelo é dada por $y|x \sim \text{Ber}(\sigma(f(x)))$ onde σ é a função sigmoide. Em resposta à falta de uma solução analítica, podemos aproximar a posteriori sobre f para qualquer conjunto de pontos de entrada usando as técnicas de inferência aproximada que vimos anteriormente. Discuta como usar a aproximação de laplace nesse caso, incluindo as fórmulas para os termos da Hessiana. Além disso, discuta como usar o resultado desse procedimento para aproximar a posteriori preditiva.

Resposta:

A probabilidade preditiva de que o evento ou a classe binária y^* seja 1, dado o conjunto de dados \bar{y} , é dada por:

$$P(y^* = 1 | \bar{y}) = \int_{\Omega} P(y^* = 1, f^* | \bar{y}) df^* \quad (12)$$

$$= \int_{\Omega} P(y^* = 1 | f^*) P(f^* | \bar{y}) df^* \quad (13)$$

$$= \int_{\Omega} \sigma(f^*) P(f^* | \bar{y}) df^* \quad (*) \quad (14)$$

$$P(f^* | \bar{y}) = \int_{\Omega} P(f^*, \bar{f} | \bar{y}) d\bar{f} \quad (15)$$

$$= \int_{\Omega} P(f^* | \bar{f}) P(\bar{f} | \bar{y}) d\bar{f} \quad (16)$$

Onde $P(f^* | \bar{f})$ é a distribuição preditiva de f^* condicional a \bar{f} . Pela questão anterior, isso é:

$$f^* | \bar{f} \sim \mathcal{N}(\mathbf{k}_{*X} (\mathbf{k}_{XX} + \sigma^2 \mathbf{I})^{-1} \bar{f}, \quad \mathbf{k}_{**} - \mathbf{k}_{*X} (\mathbf{k}_{XX} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_{X*})$$

Visto que não existe uma solução analítica para $P(\bar{f} | \bar{y})$ que é a distribuição posteriori de \bar{f} condicional a \bar{y} , usaremos a aproximação de Laplace para encontrar uma solução aproximada. Dado que a verossimilhança da classificação binária em GP's é dada pela Bernoulli: $p(y|x) \sim \text{Ber}(\sigma(f(x)))$, onde:

- y são as classes binárias.
- x são as entradas.
- σ é a função sigmoide.

Isso é feito maximizando a soma da log-verossimilhança de Bernoulli com a log-densidade da priori Gaussiana, seguindo os passos a seguir:

1. Calculamos a moda da log-posteriori, que é o ponto de máximo a posteriori (MAP) para a função latente f :

$$\mathbf{m} = \underset{\bar{f}}{\text{argmax}} \log p(\bar{y} | \bar{f}) + \log p(\bar{f})$$

A verossimilhança Bernoulli e sua log-verossimilhança são dadas por:

$$p(\bar{y}|\bar{f}) = \prod_{i=1}^N p(y_i|f_i) \quad (17)$$

$$= \prod_{i=1}^N (\sigma(f_i))^{y_i} (1 - \sigma(f_i))^{1-y_i} \quad (18)$$

$$\Rightarrow \log p(\bar{y}|\bar{f}) = \sum_{i=1}^N [y_i \log(\sigma(f_i)) + (1 - y_i) \log(1 - \sigma(f_i))] \quad (19)$$

Visto que $p(\bar{f}) = \mathcal{N}(\bar{f}|\mathbf{0}, \mathbf{K})$, temos que a log-densidade da priori é:

$$\log p(\bar{f}) = -\frac{1}{2} \bar{f}^T \mathbf{K}^{-1} \bar{f}$$

Juntando as equações acima, temos:

$$\mathbf{m} = \underset{\bar{f}}{\operatorname{argmax}} \sum_{i=1}^N [y_i \log(\sigma(f_i)) + (1 - y_i) \log(1 - \sigma(f_i))] - \frac{1}{2} \bar{f}^T \mathbf{K}^{-1} \bar{f}$$

1. Calculamos a Hessiana da log-posteriori no ponto MAP para obter a matriz de covariância da aproximação Gaussiana. A Hessiana é a matriz de segundas derivadas de uma função, e fornece uma medida da curvatura da função em torno do ponto de interesse. No nosso caso, a Hessiana da log-posteriori é dada pela soma da Hessiana da verossimilhança e da Hessiana da priori.

$$\begin{aligned} \mathbf{H} &= -\nabla \nabla \log p(\mathbf{m}|\mathbf{X}, \bar{y}) \\ &= -\nabla \nabla \left(\sum_{i=1}^N [y_i \log(\sigma(f_i)) + (1 - y_i) \log(1 - \sigma(f_i))] - \frac{1}{2} \bar{f}^T \mathbf{K}^{-1} \bar{f} \right) \end{aligned}$$

1. Com a Hessiana calculada, podemos agora determinar a matriz de covariância da aproximação Gaussiana, que é dada pelo inverso da Hessiana: \mathbf{H}^{-1}
2. Com a média \mathbf{m} e \mathbf{H}^{-1} determinados, a distribuição posteriori sobre a função latente f é aproximada como a Gaussiana:

$$p(\bar{f}|\bar{y}) \approx \mathcal{N}(\mathbf{m}, \mathbf{H}^{-1})$$

Uma vez que temos uma distribuição para f^* , a probabilidade preditiva de um novo ponto y^* dado um x^* é representada pela integral (*). Para calcular essa integral, podemos usar métodos de integração numérica, como o método de Monte Carlo. Primeiro, amostramos f^* da distribuição $p(f^*|\bar{y})$, e então calculamos a probabilidade preditiva de $y^* = 1$ como a média da função sigmoide aplicada a essas amostras.