

MAE5911/IME: Fundamentos de Estatística e Machine Learning. Prof.: Alexandre Galvão Patriota

Questão 01: Considere uma amostra aleatória $(Y_1, X_1), \dots, (Y_n, X_n)$ de (Y, X) tal que a distribuição condicional $Y | X = x \sim \mathcal{N}(\mu_\theta(x), \sigma_\theta^2(x))$, e suponha que a distribuição de X não contém informação sobre os parâmetros.

Ítem 1.1 Apresente uma **rede neural** que modele o **quantil de ordem 75%** (terceiro quartil) da distribuição condicional $Y | X = x$. O código deve ser generalizável para qualquer quantil.

Ítem 1.2 Mostre a aplicação do método nos seguintes dados simulados em R:

```
set.seed(32)

n = 1000

x = sort(runif(n, -4, 4))

y = 3 / (3 + 2|x|^3) + e^{-x^2} + cos(x) sin(x) + 0.3 \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, 1)
```

Sugestão: reescreva a esperança e variância condicionais em termos do quantil e construa a função de verossimilhança apropriada.

Resposta para Ítem 1.1

Definição da função de perda do modelo

Densidade condicional

Podemos descrever a média da Normal em função do quantil alvo $Q_q(x)$. O valor crítico da Normal padrão é dado por:

$$z_q = \Phi^{-1}(q),$$

A média condicional pode ser escrita como:

$$\mu(x) = Q_q(x) - \sigma(x) z_q.$$

Assim, o modelo probabilístico é:

$$Y | X = x \sim \mathcal{N}(\mu(x), \sigma^2(x)).$$

A densidade condicional para uma observação (x_i, y_i) é, portanto,

$$f(y_i | x_i; Q_q, \sigma) = \frac{1}{\sqrt{2\pi} \sigma(x_i)} \exp \left\{ -\frac{(y_i - Q_q(x_i) + \sigma(x_i) z_q)^2}{2\sigma^2(x_i)} \right\}.$$

Log-verossimilhança

Tomando log da expressão acima, obtemos a log-verossimilhança para uma única observação:

$$\ell_i(Q_q, \sigma) = \log f(y_i | x_i; Q_q, \sigma) = -\frac{1}{2} \log(2\pi) - \log \sigma(x_i) - \frac{(y_i - Q_q(x_i) + \sigma(x_i)z_q)^2}{2\sigma^2(x_i)}.$$

Somando sobre todas as observações, a log-verossimilhança total é:

$$\ell(Q_q, \sigma) = \sum_{i=1}^n \ell_i(Q_q, \sigma) = -\frac{n}{2} \log(2\pi) - \sum_{i=1}^n \log \sigma(x_i) - \sum_{i=1}^n \frac{(y_i - Q_q(x_i) + \sigma(x_i)z_q)^2}{2\sigma^2(x_i)}.$$

O negativo da log-verossimilhança (função de perda) é, portanto,

$$\mathcal{L}_{\text{NLL}}(Q_q, \sigma) = -\ell(Q_q, \sigma) = \frac{n}{2} \log(2\pi) + \sum_{i=1}^n \log \sigma(x_i) + \sum_{i=1}^n \frac{(y_i - Q_q(x_i) + \sigma(x_i)z_q)^2}{2\sigma^2(x_i)}.$$

O termo

$$\frac{n}{2} \log(2\pi)$$

não depende de Q_q nem de $\sigma(x)$, e portanto não altera o ponto de mínimo nem os gradientes no treinamento da rede neural. Assim, podemos descartá-lo e usar a forma simplificada da log-verossimilhança:

$$\ell_i(Q_q, \sigma) = -\log \sigma(x_i) - \frac{(y_i - Q_q(x_i) + \sigma(x_i)z_q)^2}{2\sigma^2(x_i)}.$$

Portanto, a função de perda que iremos usar no treinamento da rede, negativo da log-verossimilhança, é:

$$\mathcal{L}_{\text{NLL}}(Q_q, \sigma) = \sum_{i=1}^n \left[\log \sigma(x_i) + \frac{(y_i - Q_q(x_i) + \sigma(x_i)z_q)^2}{2\sigma^2(x_i)} \right].$$

Definição da Arquitetura

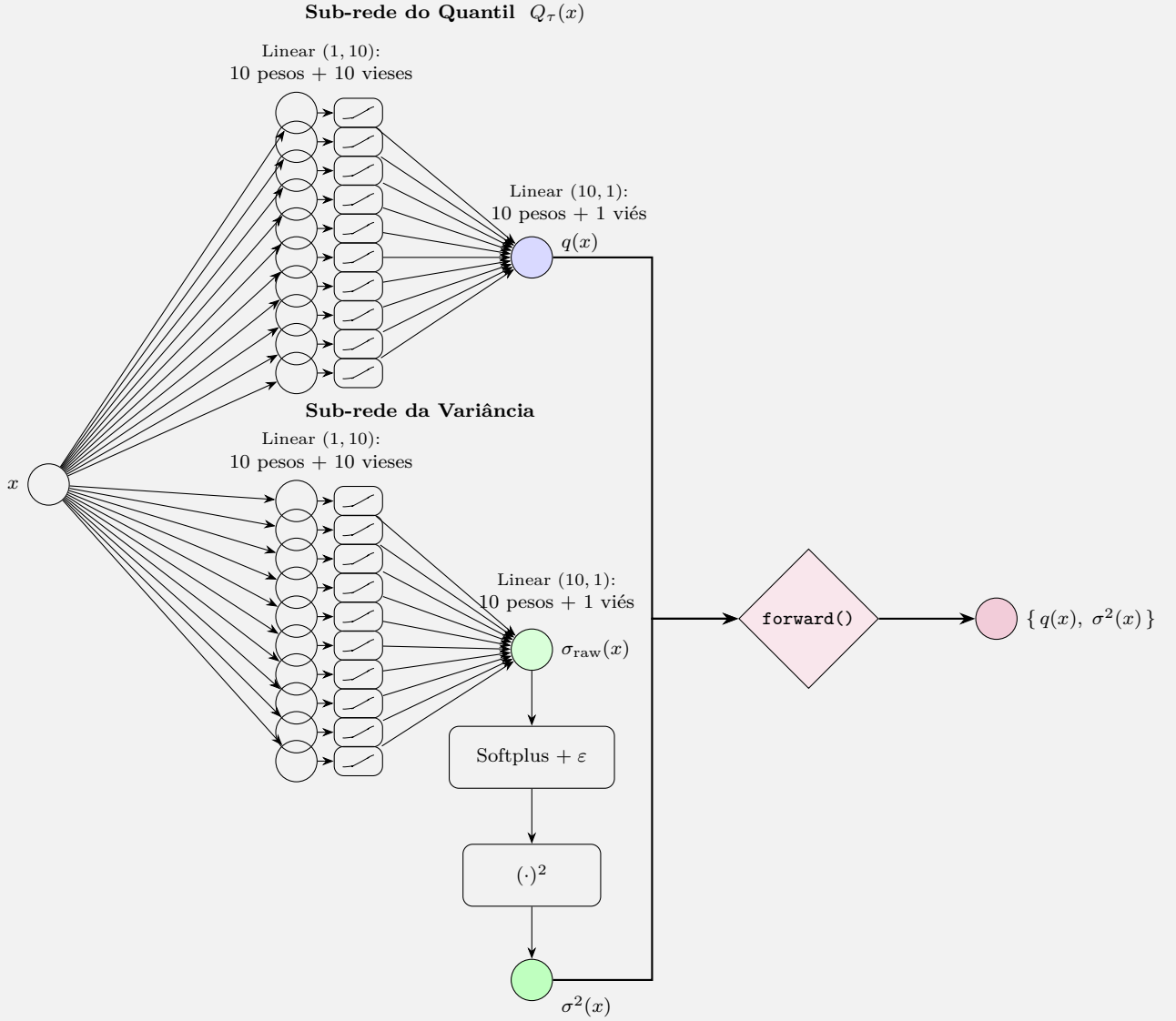


Figura 1: Arquitetura da rede neural com 10 neurônios ocultos em cada subrede.

Parâmetro	Dimensão	Nº de parâmetros
q_net.0.weight	10×1	10
q_net.0.bias	10	10
q_net.2.weight	1×10	10
q_net.2.bias	1	1
sigma_net.0.weight	10×1	10
sigma_net.0.bias	10	10
sigma_net.2.weight	1×10	10
sigma_net.2.bias	1	1
Total	—	62

Tabela 1: Parâmetros treináveis da rede neural utilizada no experimento.

```

1 library(torch)
2
3 QuantileNormalNet <- nn_module(
4   "QuantileNormalNet",
5
6   initialize = function(input_dim,
7                           hidden_q    = config$hidden_mu,
8                           hidden_sigma = config$hidden_sigma) {
9

```

```

10 # Rede do quantil Q_q(x)
11 self$q_net <- nn_sequential(
12   nn_linear(input_dim, hidden_q),
13   nn_gelu(),
14   nn_linear(hidden_q, 1)
15 )
16
17 # Rede da variância sigma(x)
18 self$sigma_net <- nn_sequential(
19   nn_linear(input_dim, hidden_sigma),
20   nn_gelu(),
21   nn_linear(hidden_sigma, 1)
22 )
23 },
24
25 forward = function(x) {
26   # x: tensor [batch_size, 1]
27
28   # Predição do quantil Q_q(x)
29   q_pred <- self$q_net(x)
30
31   # Predição da variância sigma(x)
32   sigma_raw <- self$sigma_net(x)
33   sigma <- nnf_softplus(sigma_raw) + 1e-4
34   sigma2 <- sigma$pow(2)
35
36   list(
37     q = q_pred,
38     sigma2 = sigma2
39   )
40 }
41 )

```

Listing 1: Estrutura do modelo QuantileNormalNet em torch para R

```

1 nll_normal_from_quantile <- function(Qq, sigma2, y, q_level, eps = 1e-6) {
2
3   Qq <- torch_tensor(Qq, dtype = torch_float())
4   sigma2 <- torch_tensor(sigma2, dtype = torch_float())
5   y <- torch_tensor(y, dtype = torch_float())
6
7   sigma <- sigma2$sqrt()
8
9   z_q <- qnorm(q_level)
10  z_q_t <- torch_tensor(z_q, dtype = torch_float(), device = Qq$device)
11
12  mu <- Qq - sigma * z_q_t
13
14  term1 <- torch_log(sigma)
15  term2 <- (y - mu)$pow(2) / (2 * sigma$pow(2))
16
17  (term1 + term2)$mean()
18 }

```

Listing 2: Função de perda NLL para o modelo Normal parametrizado pelo quantil $Q_q(x)$

Resposta para Ítem 1.2

Treinando o modelo proposto na Seção anterior com os dados do enunciado, obtemos os seguintes resultados:

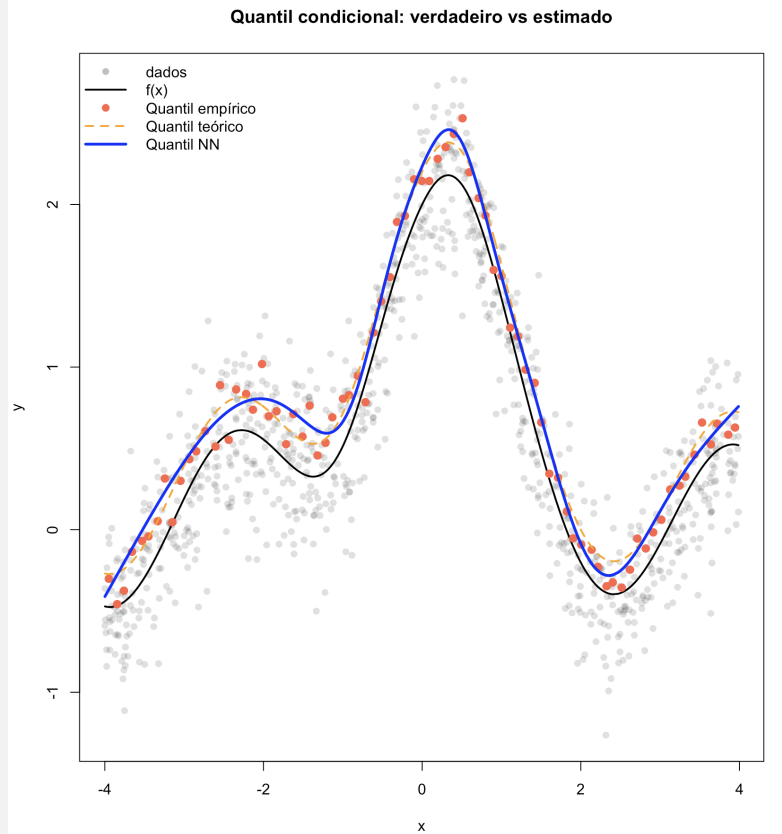


Figura 2: Quantil condicional verdadeiro vs. estimado.

A Figura 2 apresenta a comparação entre o quantil condicional verdadeiro, o quantil empírico observado nos dados e o quantil estimado pela rede neural. A curva aprendida (em azul) acompanha o comportamento oscilatório do quantil teórico. Isso evidencia que a rede foi capaz de aprender a dependência não linear de $Q_\tau(x)$.

A evolução da função de perda (NLL), cuja trajetória está ilustrada na Figura 3 mostra que as curvas de *train* e *test* permanecem próximas ao longo de todas as épocas, indicando ausência de sobreajuste e capacidade de generalização.

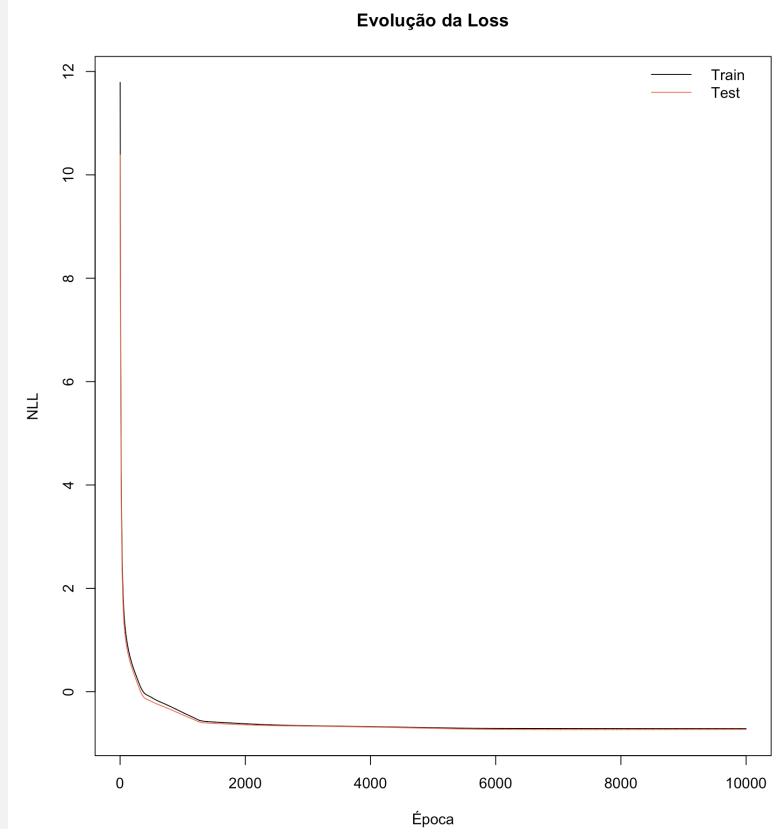


Figura 3: Evolução da NLL durante o treinamento.

```

1 source("Config.R")
2 library(torch)
3
4 set.seed(config$seed)
5
6 n <- config$n
7 x <- sort(runif(n, config$x_min, config$x_max))
8
9 f_x <- 3/(3 + 2*abs(x)^3) +
10   exp(-x^2) +
11   cos(x)*sin(x)
12
13 sigma_true <- rep(config$noise_sd, length(x))
14 sigma2_true <- rep(config$noise_sd^2, length(x))
15
16 y <- f_x + rnorm(n, mean = 0, sd = sigma_true)
17
18 X_full <- torch_tensor(matrix(x, ncol = 1), dtype = torch_float())
19 Y_full <- torch_tensor(matrix(y, ncol = 1), dtype = torch_float())
20
21 n_train <- round(config$p_train * n)
22 idx <- sample(1:n)
23
24 idx_train <- idx[1:n_train]
25 idx_test <- idx[(n_train + 1):n]
26
27 x_train <- x[idx_train]
28 y_train <- y[idx_train]
29
30 x_test <- x[idx_test]
31 y_test <- y[idx_test]
32
33 X_train <- torch_tensor(matrix(x_train, ncol = 1), dtype = torch_float())
34 Y_train <- torch_tensor(matrix(y_train, ncol = 1), dtype = torch_float())
35
36 X_test <- torch_tensor(matrix(x_test, ncol = 1), dtype = torch_float())
37 Y_test <- torch_tensor(matrix(y_test, dtype = torch_float()))

```

Listing 3: Geração dos dados para o Item 1.2: modelo homoscedástico $y = f(x) + \sigma\epsilon$.

```

1 library(torch)
2
3 treinar_normal_quantil <- function(model,
4                                   q_level = config$q,
5                                   epochs = config$num_epochs,
6                                   lr      = config$lr) {
7
8   print_every <- config$print_every
9   plot_every  <- config$plot_every
10  num_bins    <- 80
11
12  optimizer <- optim_adamw(
13    model$parameters,
14    lr = lr,
15    weight_decay = 1e-3
16  )
17
18  loss_history_train <- numeric(epochs)
19  loss_history_test  <- numeric(epochs)
20
21  par(mfrow = c(1, 2))
22
23  for (epoch in 1:epochs) {
24
25    model$train()
26    optimizer$zero_grad()
27
28    out_train <- model(X_train)
29    Qq_train  <- out_train$q
30    sigma2_tr <- out_train$sigma2
31
32    loss_train <- nll_normal_from_quantile(
33      Qq      = Qq_train,
34      sigma2  = sigma2_tr,
35      y       = Y_train,
36      q_level = q_level
37    )
38
39    loss_train$backward()
40    optimizer$step()
41
42    model$eval()
43    out_test <- model(X_test)
44
45    Qq_test  <- out_test$q
46    sigma2_te <- out_test$sigma2
47
48    loss_test <- nll_normal_from_quantile(
49      Qq      = Qq_test,
50      sigma2  = sigma2_te,
51      y       = Y_test,
52      q_level = q_level
53    )
54
55    loss_history_train[epoch] <- loss_train$item()
56    loss_history_test[epoch]  <- loss_test$item()
57
58    if (epoch %% print_every == 0) {
59      cat("Época:", epoch,
60          "\n  - NLL train:", loss_train$item(),
61          "\n  - NLL test:", loss_test$item(), "\n")
62    }
63
64    if (epoch %% plot_every == 0 || epoch == epochs) {
65

```

```

66 out_full <- model(X_full)
67 Qq_full <- as.numeric(out_full$q$squeeze())
68 sigma2_hat <- as.numeric(out_full$sigma2$squeeze())
69 sigma_hat <- sqrt(sigma2_hat)
70
71 z_q <- qnorm(q_level)
72 q_theo <- f_x + z_q * sigma_true # sigma_true é constante no item 1.2
73
74 ylim_all <- range(y, f_x, q_theo, Qq_full)
75
76 plot(x, y,
77      pch = 16,
78      col = rgb(0,0,0,0.15),
79      ylim = ylim_all,
80      xlab = "x", ylab = "y",
81      main = "Quantil condicional: verdadeiro vs estimado")
82
83 lines(x, f_x, col = "black", lwd = 2)
84
85 plot_conditional_quantile(
86   x, y, q_level,
87   num_bins = num_bins,
88   col = "tomato", pch = 19, cex = 1.1
89 )
90
91 lines(x, q_theo, col = "orange", lwd = 2, lty = 2)
92 lines(x, Qq_full, col = "blue", lwd = 3)
93
94 legend("topleft",
95       legend = c("dados", "f(x)", "Quantil empírico",
96                 "Quantil teórico", "Quantil NN"),
97       col = c(rgb(0,0,0,0.3), "black", "tomato", "orange", "blue"),
98       lwd = c(NA, 2, NA, 2, 3),
99       lty = c(NA, 1, NA, 2, 1),
100      pch = c(16, NA, 19, NA, NA),
101      bty = "n")
102
103 plot(1:epoch, loss_history_train[1:epoch], type = "l",
104      xlab = "Época", ylab = "NLL",
105      ylim = range(c(loss_history_train[1:epoch],
106                    loss_history_test[1:epoch])),
107      main = "Evolução da Loss")
108
109 lines(1:epoch, loss_history_test[1:epoch], col = "tomato")
110
111 legend("topright",
112       legend = c("Train", "Test"),
113       col = c("black", "tomato"),
114       lty = 1, bty = "n")
115
116 Sys.sleep(0.05)
117 }
118 }
119
120 invisible(list(
121   model = model,
122   loss_train = loss_history_train,
123   loss_test = loss_history_test
124 ))
125 }
126
127 plot_conditional_quantile <- function(x, y, q, num_bins = 80,
128                                       col = "darkgreen", pch = 19, cex = 1.1) {
129
130   breaks <- seq(min(x), max(x), length.out = num_bins)
131   qx <- c(); qy <- c()
132
133   for (i in 1:(length(breaks) - 1)) {

```



```

134     idx <- which(x >= breaks[i] & x < breaks[i + 1])
135     if (length(idx) > 0) {
136       qx <- c(qx, mean(x[idx]))
137       qy <- c(qy, quantile(y[idx], q))
138     }
139   }
140
141   points(qx, qy, col = col, pch = pch, cex = cex)
142   invisible(list(x = qx, y = qy))
143 }

```

Listing 4: Função de treinamento para o modelo Normal-Quantil no caso homoscedástico (Item 1.2).

```

1  config <- list(
2
3  # =====
4  # Quantil alvo
5  # =====
6  q = 0.75,
7
8  # =====
9  # Hiperparâmetros de treino
10 # =====
11 num_epochs = 10000,
12 lr          = 1e-3,
13 print_every = 10L,
14 plot_every  = 10L,
15
16 # =====
17 # Arquitetura da rede
18 # =====
19 input_dim   = 1L,
20 hidden_mu   = 10L, # 10 neurônios para sub-rede do quantil
21 hidden_sigma = 10L, # 10 neurônios para sub-rede da variância
22
23 # =====
24 # Configuração dos dados
25 # =====
26 seed        = 32,
27 n            = 1000,
28 x_min        = -4,
29 x_max        = 4,
30 p_train      = 0.7,
31
32 # =====
33 # Variância verdadeira usada na simulação
34 # =====
35 noise_sd = 0.3,
36 )

```

Listing 5: Arquivo Config.R utilizado no Item 1.2.

Os códigos deste estudo estão disponibilizados em <http://academic-codex.github.io/MAE5911-Estatistica-e-Machine-Learning>.

