

# Métodos geométricos aplicados a ciências da informação

Fábio C. C. Meneghetti

Instituto de Matemática, Estatística e Computação Científica  
Universidade Estadual de Campinas

22 de fevereiro de 2022

# Parte I

## Motivacao

# História

- Claude Shannon (1916–2001): entropia como medida de informação de uma variável aleatória:  $H(X) = -\sum_i p_i \log p_i$ .

# História

- Claude Shannon (1916–2001): entropia como medida de informação de uma variável aleatória:  $H(X) = -\sum_i p_i \log p_i$ .
- Fisher (1890–1962): medida de informação sobre um parâmetro  $\theta$  que uma variável aleatória carrega:  $I(\theta) = \int_{\mathcal{R}} p(x) \log \frac{dp}{d\theta}(x)^2 dx$

# História

- Claude Shannon (1916–2001): entropia como medida de informação de uma variável aleatória:  $H(X) = -\sum_i p_i \log p_i$ .
- Fisher (1890–1962): medida de informação sobre um parâmetro  $\theta$  que uma variável aleatória carrega:  $I(\theta) = \int_{\mathbb{R}} p(x) \log \frac{dp}{d\theta}(x)^2 dx$
- C. R. Rao (1920–): a informação de Fisher sobre múltiplos parâmetros  $(\theta_1, \dots, \theta_d)$  é uma métrica Riemanniana!

# História

- Claude Shannon (1916–2001): entropia como medida de informação de uma variável aleatória:  $H(X) = -\sum_i p_i \log p_i$ .
- Fisher (1890–1962): medida de informação sobre um parâmetro  $\theta$  que uma variável aleatória carrega:  $I(\theta) = \int_{\mathcal{X}} p(x) \log \frac{dp}{d\theta}(x)^2 dx$
- C. R. Rao (1920–): a informação de Fisher sobre múltiplos parâmetros  $(\theta_1, \dots, \theta_d)$  é uma métrica Riemanniana!
  - *metrica da informacao de Fisher*

# História

- Claude Shannon (1916–2001): entropia como medida de informação de uma variável aleatória:  $H(X) = -\sum_i p_i \log p_i$ .
- Fisher (1890–1962): medida de informação sobre um parâmetro  $\theta$  que uma variável aleatória carrega:  $I(\theta) = \int_{\mathcal{X}} p(x) \log \frac{dp}{d\theta}(x)^2 dx$
- C. R. Rao (1920–): a informação de Fisher sobre múltiplos parâmetros  $(\theta_1, \dots, \theta_d)$  é uma métrica Riemanniana!
  - *metrica da informacao de Fisher*
  - ela induz uma geometria sobre as distribuições de probabilidade

# História

- Claude Shannon (1916–2001): entropia como medida de informação de uma variável aleatória:  $H(X) = -\sum_i p_i \log p_i$ .
- Fisher (1890–1962): medida de informação sobre um parâmetro  $\theta$  que uma variável aleatória carrega:  $I(\theta) = \int_{\mathcal{X}} p(x) \log \frac{dp}{d\theta}(x)^2 dx$
- C. R. Rao (1920–): a informação de Fisher sobre múltiplos parâmetros  $(\theta_1, \dots, \theta_d)$  é uma métrica Riemanniana!
  - *metrica da informacao de Fisher*
  - ela induz uma geometria sobre as distribuições de probabilidade
  - isso nos permite falar sobre distâncias e curvaturas no espaço das distribuições



# História

- Claude Shannon (1916–2001): entropia como medida de informação de uma variável aleatória:  $H(X) = -\sum_i p_i \log p_i$ .
- Fisher (1890–1962): medida de informação sobre um parâmetro  $\theta$  que uma variável aleatória carrega:  $I(\theta) = \int_{\mathcal{X}} p(x) \log \frac{dp}{d\theta}(x)^2 dx$
- C. R. Rao (1920–): a informação de Fisher sobre múltiplos parâmetros  $(\theta_1, \dots, \theta_d)$  é uma métrica Riemanniana!
  - *metrica da informacao de Fisher*
  - ela induz uma geometria sobre as distribuições de probabilidade
  - isso nos permite falar sobre distâncias e curvaturas no espaço das distribuições
  - área de pesquisa: *geometria da informacao*

## Por exemplo:

- no caso das distribuições gaussianas univariadas...

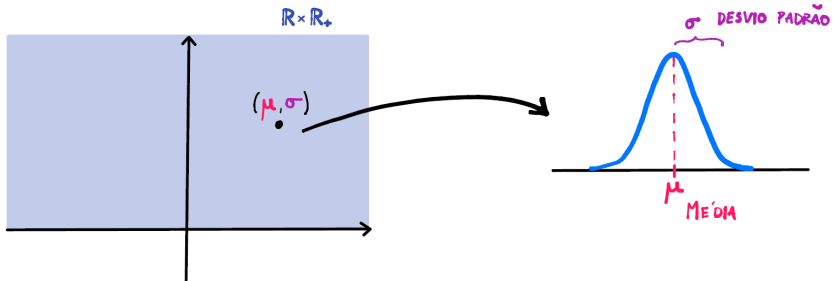
$$p(x; \mu, \sigma^2) = \frac{1}{\sigma \sqrt{2\pi}} \exp \left( -\frac{(x - \mu)^2}{2\sigma^2} \right);$$

# Por exemplo:

- no caso das distribuições gaussianas univariadas...

$$p(x; \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right);$$

o espaço de parâmetros é



- com a métrica da informação de Fisher, obtemos uma geometria hiperbólica! (versão deformada do meio-plano de Poincaré)

- com a métrica da informação de Fisher, obtemos uma geometria hiperbólica! (versão deformada do meio-plano de Poincaré)



## Parte II

# Teoria

# Distribuições de probabilidade

distribuições de probabilidade

# Distribuições de probabilidade

distribuições de probabilidade

discretas:  $X \in \{x_1, \dots, x_n\}$ ,  $P[X = x_i] = p_i$ ,  $\sum_i p_i = 1$ .



# Distribuições de probabilidade

distribuições de probabilidade

discretas:  $X \in \{x_1, \dots, x_n\}$ ,  $P[X = x_i] = p_i$ ,  $\sum_i p_i = 1$ .

contínuas:  $X \in \mathbb{R}$ ,  $P[X \in A] = \int_A p(x) dx$ .

# Distribuições de probabilidade

distribuições de probabilidade

discretas:  $X \in \{x_1, \dots, x_n\}$ ,  $P[X = x_i] = p_i$ ,  $\sum_i p_i = 1$ .

contínuas:  $X \in \mathbb{R}$ ,  $P[X \in A] = \int_A p(x) dx$ .

de forma mais geral, temos um espaço de probabilidade  $(\mathbb{R}; P)$ , e uma medida - nita dominante .

# Distribuições de probabilidade

distribuições de probabilidade

discretas:  $X \in \{x_1, \dots, x_n\}$ ,  $P[X = x_i] = p_i$ ,  $\sum_i p_i = 1$ .

contínuas:  $X \in \mathbb{R}$ ,  $P[X \in A] = \int_A p(x) dx$ .

de forma mais geral, temos um espaço de probabilidade  $(\mathbb{R}; P)$ , e uma medida - nita dominante  $\mu$ .

A função densidade é a derivada de Radon-Nikodym  $p(x) = \frac{dP}{d\mu}(x)$ ,  
 $p: \mathbb{R} \rightarrow \mathbb{R}_+$  que satisfaz  $P(A) = \int_A p(x) d\mu(x)$ .

# Distribuições de probabilidade

distribuições de probabilidade

discretas:  $X \in \{x_1, \dots, x_n\}$ ,  $P[X = x_i] = p_i$ ,  $\sum_i p_i = 1$ .

contínuas:  $X \in \mathbb{R}$ ,  $P[X \in A] = \int_A p(x) dx$ .

de forma mais geral, temos um espaço de probabilidade  $(\mathbb{R}; P)$ , e uma medida - nita dominante  $\mu$ .

A função densidade é a derivada de Radon-Nikodym  $p(x) = \frac{dP}{d\mu}(x)$ ,  
 $p: \mathbb{R} \rightarrow \mathbb{R}_+$  que satisfaz  $P(A) = \int_A p(x) d\mu(x)$ .

$X$  enumerável e  $\mu$  medida de contagem  $\Rightarrow$  distribuição discreta,  
 $p =$  função massa,  $\int_{x \in A} p(x) d\mu(x) = \sum_{x \in A} p(x)$

# Distribuições de probabilidade

distribuições de probabilidade

discretas:  $X \in \{x_1, \dots, x_n\}$ ,  $P[X = x_i] = p_i$ ,  $\sum_i p_i = 1$ .

contínuas:  $X \in \mathbb{R}$ ,  $P[X \in A] = \int_A p(x) dx$ .

de forma mais geral, temos um espaço de probabilidade  $(\mathbb{R}; P)$ , e uma medida  $\mu$  - nita dominante.

A função densidade é a derivada de Radon-Nikodym  $p(x) = \frac{dP}{d\mu}(x)$ ,  
 $p: \mathbb{R} \rightarrow \mathbb{R}_+$  que satisfaz  $P(A) = \int_A p(x) d\mu(x)$ .

$X$  enumerável e  $\mu$  medida de contagem) = distribuição discreta,  
 $p =$  função massa,  $\int_{x \in A} p(x) d\mu(x) = \sum_{x \in A} p(x)$

$X \in \mathbb{R}^n$  e  $\mu$  medida de Lebesgue) = distribuição contínua,  $p =$   
 função densidade

# Modelos estatísticos

um modelo estatístico é uma família parametrizada de distribuições  $P = f(\theta)$  no espaço  $(X; \mathcal{F})$ .

# Modelos estatísticos

um modelo estatístico é uma família parametrizada de distribuições  $P = f(\theta)$  no espaço  $(X; F)$ .

$\Theta \subset \mathbb{R}^d$  conjunto aberto de parâmetros.

# Modelos estatísticos

um modelo estatístico é uma família parametrizada de distribuições  $P = f(\theta)$  no espaço  $(X; \mathcal{F})$ .

$\mathbb{R}^d$  conjunto aberto de parâmetros.

na vida real costuma-se tomar como modelo a família de funções de densidade  $p = \frac{dP}{d\theta}$ .



# Modelos estatísticos

um modelo estatístico é uma família parametrizada de distribuições  $P = f(\theta)$  no espaço  $(X; F)$ .

$\mathbb{R}^d$  conjunto aberto de parâmetros.

na vida real costuma-se tomar como modelo a família de funções de densidade  $p = \frac{dP}{d\theta}$ .

vamos considerar modelos estatísticos regulares

# Modelos estatísticos

um modelo estatístico é uma família parametrizada de distribuições  $P = f(\theta)$  no espaço  $(X; \mathcal{F})$ .

$\Theta \subset \mathbb{R}^d$  conjunto aberto de parâmetros.

na vida real costuma-se tomar como modelo a família de funções de densidade  $p = \frac{dP}{d\mu}$ .

vamos considerar modelos estatísticos regulares  
 a função  $p$  é  $C^1$

# Modelos estatísticos

um modelo estatístico é uma família parametrizada de distribuições  $P = f(\theta) : \theta \in \Theta$  no espaço  $(X; \mathcal{F})$ .

$\Theta \subset \mathbb{R}^d$  conjunto aberto de parâmetros.

na vida real costuma-se tomar como modelo a família de funções de densidade  $p = \frac{dP}{d\mu} : \theta \in \Theta$ .

vamos considerar modelos estatísticos regulares

- a função  $p$  é  $C^1$
- $p(x) > 0$  para todo  $x$ ;

# Métrica Riemanniana

Vamos agora adicionar uma estrutura Riemanniana ao espaço  $\mathbb{P}$  .

# Métrica Riemanniana

Vamos agora adicionar uma estrutura Riemanniana ao espaço  $\mathcal{P}$ .  
 denote por  $\eta(x) := \log p(x)$  a função log-probabilidade

# Métrica Riemanniana

Vamos agora adicionar uma estrutura Riemanniana ao espaço  $\mathcal{P}$ .

denote por  $\eta(x) := \log p(x)$  a função log-probabilidade  
 a métrica de Fisher a métrica Riemanniana (produto interno em  $T_x$ ) dada por

$$g(V; W) := \int_{\mathcal{X}} p(x) \frac{\partial \eta(x)}{\partial V} \frac{\partial \eta(x)}{\partial W} dx$$

# Métrica Riemanniana

Vamos agora adicionar uma estrutura Riemanniana ao espaço  $\mathcal{P}$ .

denote por  $\eta(x) := \log p(x)$  a função log-probabilidade  
 a métrica de Fisher a métrica Riemanniana (produto interno em  $T_x$ ) dada por

$$g(V; W) := \int_{\mathcal{X}} p(x) \frac{\partial \eta(x)}{\partial \theta_i} \frac{\partial \eta(x)}{\partial \theta_j} d(x)$$

na base coordenada local  $\{e_i\}_i$ , a matriz da métrica é chamada de matriz de Fisher  $I(\theta)$ , com elementos

$$g_{ij}(\theta) := g(e_i; e_j) = \int_{\mathcal{X}} p(x) \frac{\partial \eta(x)}{\partial \theta_i} \frac{\partial \eta(x)}{\partial \theta_j} d(x)$$

# Métrica Riemanniana

Vamos agora adicionar uma estrutura Riemanniana ao espaço  $\mathcal{P}$ .

denote por  $\eta(x) := \log p(x)$  a função log-probabilidade  
 a métrica de Fisher a métrica Riemanniana (produto interno em  $T_x$ ) dada por

$$g(V; W) := \int_{\mathcal{X}} p(x) \frac{\partial \eta(x)}{\partial \theta} \frac{\partial \eta(x)}{\partial \theta} d(x)$$

na base coordenada local  $\{e_i\}_i$ , a matriz da métrica é chamada de matriz de Fisher  $I(\theta)$ , com elementos

$$g_{ij}(\theta) := g(e_i; e_j) = \int_{\mathcal{X}} p(x) \frac{\partial \eta(x)}{\partial \theta_i} \frac{\partial \eta(x)}{\partial \theta_j} d(x)$$

$I(\theta)$  é simétrica e positiva-definita.



# Interpretações da métrica de Fisher

Em estatística

# Interpretações da métrica de Fisher

Em estatística

O escorede  $\eta$  em  $x$  e o gradiente da função log-probabilidade:

$$s(x) = \eta \quad \text{log} p(x) = \left( \frac{\partial}{\partial \eta}; \dots; \frac{\partial}{\partial \eta_d} \right)^T(x)$$

mede sensibilidade a mudanças nos parâmetros da função log-probabilidade

# Interpretações da métrica de Fisher

## Em estatística

O escore de  $x$  é o gradiente da função log-probabilidade:

$$s(x) = \left( \frac{\partial}{\partial \theta_1}, \dots, \frac{\partial}{\partial \theta_d} \right)^T \log p(x)$$

mede sensibilidade a mudanças nos parâmetros da função log-probabilidade

A matriz de Fisher é a matriz de covariância do escore:

$$I(\theta) = \text{cov}(s; s) = E[s s^T]$$

e um limitante inferior para a covariância de um estimador não-enviesado  $\hat{\theta}$  (Limitante de Cramér-Rao):

$$\text{cov}(\hat{\theta}) \succeq I(\theta)^{-1}$$

em teoria da informação

em teoria da informação

entropia relativa (divergência de Kullback-Leibler):

$$D_{\text{KL}}(p \parallel q) = \int_{\mathcal{X}} p(x) \log \frac{p(x)}{q(x)} d(x)$$

## em teoria da informação

entropia relativa (divergência de Kullback-Leibler):

$$D_{KL}(p \parallel q) = \int_{\mathcal{X}} p(x) \log \frac{p(x)}{q(x)} d(x)$$

diz quantos bits (na verdade nats), em média, são necessários para codificar  $p$  usando um código otimizado para codificar  $q$

- em teoria da informação

- entropia relativa (divergência de Kullback-Leibler):

$$D_{\text{KL}}(p \parallel q) = \int_X p(x) \log \frac{p(x)}{q(x)} d(x)$$

- diz quantos bits (na verdade nats), em media, são necessários para codificar  $p$  usando um código otimizado para codificar  $q$
- a matriz de Fisher é o Hessiano diagonal da entropia relativa:

$$g_{ij}(\theta) = \frac{\partial^2}{\partial \theta_i \partial \theta_j} D_{\text{KL}}(p_{\theta} \parallel p)$$

# Unicidade

uma estatística e um mapa mensurável :  $X \rightarrow Y$



# Unicidade

uma estatística e um mapa mensurável  $\mu : X \rightarrow Y$   
 induz uma família de distribuições de probabilidade empurradas  $\mathcal{P}$   
 em  $Y$ : 
$$P(A) := P(\mu^{-1}(A))$$

# Unicidade

uma estatística e um mapa mensurável  $\gamma : X \rightarrow Y$   
 induz uma família de distribuições de probabilidade empurradas  $\mathcal{P}$   
 em  $Y$ :  $P(A) := P(\gamma^{-1}(A))$   
 por sua vez, obtemos novas funções densidade  $\frac{dP}{d\mu}$

# Unicidade

uma estatística e um mapa mensurável  $\tau : X \rightarrow Y$

induz uma família de distribuições de probabilidade empurradas  $\mathcal{P}$  em  $Y$ :  $P(A) := P(\tau^{-1}(A))$

por sua vez, obtemos novas funções densidade  $\frac{dP}{d\mu}$

a estatística é dita suficiente se  $p(x) = \eta(\tau(x))h(x)$  para alguma função  $h$  independente de  $\eta$ .

# Unicidade

uma estatística e um mapa mensurável  $\mu : X \rightarrow Y$

induz uma família de distribuições de probabilidade empurradas  $\mathcal{P}$  em  $Y$ :  $P(A) := P(\mu^{-1}(A))$

por sua vez, obtemos novas funções densidade  $\frac{d}{d\mu} P$

a estatística é dita suficiente se  $p(x) = p(\mu(x))h(x)$  para alguma função  $h$  independente de  $\mu$ .

A métrica de Fisher é a única métrica Riemanniana, a menos de uma constante, invariante por estatísticas suficientes.

# Geodésicas

geodésicas são "linhas retas" nas variedades Riemannianas

# Geodésicas

geodésicas são "linhas retas" nas variedades Riemannianas  
dados dois pontos  $p, p_0$ , a curva  $\gamma : [0; 1] \rightarrow P$  ligando-os, que  
minimiza comprimento, é um segmento de geodésica

# Geodésicas

geodésicas são "linhas retas" nas variedades Riemannianas  
dados dois pontos  $p, p_0$ , a curva  $\gamma : [0; 1] \rightarrow P$  ligando-os, que  
minimiza comprimento, é um segmento de geodésica

o comprimento é dado por 
$$L(\gamma) = \int_0^1 \sqrt{g(\dot{\gamma}(t); \dot{\gamma}(t))} dt$$

# Geodésicas

geodésicas são "linhas retas" nas variedades Riemannianas  
dados dois pontos  $p ; p_0$ , a curva  $\gamma : [0; 1] \rightarrow P$  ligando-os, que  
minimiza comprimento, é um segmento de geodésica

o comprimento é dado por  $l(\gamma) = \int_0^1 \sqrt{g_{(t)}(\dot{\gamma}(t); \dot{\gamma}(t))} dt$   
essa noção de distância geodésica é chamada distância de  
Fisher-Rao<sup>1</sup> na geometria da informação:

$$d_{FR}(p ; p^0) = \min_{\gamma} l(\gamma) : \gamma(0) = p ; \gamma(1) = p^0$$

<sup>1</sup>quando a variedade é completa e conexa por caminhos



# Geodésicas

geodésicas são "linhas retas" nas variedades Riemannianas dados dois pontos  $p ; p_0$ , a curva  $\gamma : [0; 1] \rightarrow P$  ligando-os, que minimiza comprimento, é um segmento de geodésica

o comprimento é dado por 
$$L(\gamma) = \int_0^1 \sqrt{g_{(t)}(\dot{\gamma}(t); \dot{\gamma}(t))} dt$$
 essa noção de distância geodésica chamada distância de Fisher-Rao<sup>1</sup> na geometria da informação:

$$d_{FR}(p ; p^0) = \min_{\gamma} L(\gamma) : \gamma(0) = p ; \gamma(1) = p^0$$

formalmente, são curvas que têm derivada covariante zero:  $\nabla_{\dot{\gamma}} \dot{\gamma} = 0$

<sup>1</sup>quando a variedade é completa e conexa por caminhos

# Geodésicas

geodésicas são "linhas retas" nas variedades Riemannianas dados dois pontos  $p$  ;  $p^0$ , a curva  $\gamma : [0; 1] \rightarrow P$  ligando-os, que minimiza comprimento, é um segmento de geodésica

o comprimento é dado por 
$$L(\gamma) = \int_0^1 \sqrt{g_{ij}(\gamma(t)) \dot{\gamma}^i(t) \dot{\gamma}^j(t)} dt$$
 essa noção de distância geodésica é chamada distância de Fisher-Rao<sup>1</sup> na geometria da informação:

$$d_{FR}(p ; p^0) = \min_{\gamma} L(\gamma) : \gamma(0) = p ; \gamma(1) = p^0$$

formalmente, são curvas que têm derivada covariante zero:  $\nabla_{\dot{\gamma}} \dot{\gamma} = 0$

onde  $\nabla$  é a conexão dada pelos símbolos de Christoffel

$$\nabla_{\dot{\gamma}} \dot{\gamma} = \dot{\gamma}^i \left( \dot{\gamma}^j \frac{\partial}{\partial x^i} + \frac{1}{2} \dot{\gamma}^i \dot{\gamma}^j \Gamma_{ij}^k \right) \dot{\gamma}^k$$

<sup>1</sup>quando a variedade é completa e conexa por caminhos

# Famílias exponenciais

uma família exponencial  $p : \mathbb{R}^d \rightarrow \mathbb{R}^d$  com parâmetros naturais  
 $\mathbb{R}^d$  e dada por

$$p(x) := \exp \{ \eta(x); \theta \in \mathbb{R}^d \} + k(x); \quad x \in X$$

# Famílias exponenciais

uma família exponencial  $p : \mathbb{R}^d \rightarrow \mathbb{R}^+$  com parâmetros naturais  $\eta \in \mathbb{R}^d$  e dada por

$$p(x) := \exp(\eta^T t(x)); \quad \int_{\mathcal{X}} \exp(\eta^T t(x)) h(x) dx = F(\eta) + k(\eta); \quad x \in \mathcal{X}$$

$t(x)$  estatística suficiente

# Famílias exponenciais

uma família exponencial  $p : \mathbb{R}^d \rightarrow \mathbb{R}^d$  com parâmetros naturais  
 $\mathbb{R}^d$  e dada por

$$p(x) := \exp \{ \eta(x); \theta \in \mathbb{R}^d \} + k(x); \quad x \in X$$

$\eta(x)$  estatística suficiente

$F(\theta)$  função estritamente convexa

# Famílias exponenciais

uma família exponencial  $\mathcal{P} \subset \mathcal{G}$  com parâmetros naturais  
 $\theta \in \mathbb{R}^d$  e dada por

$$p(x) := \exp \{ \eta(\theta) t(x) \}; \quad \eta(\theta) = F(\theta) + k(x); \quad x \in X$$

$t(x)$  estatística suficiente

$F(\cdot)$  função estritamente convexa

$k(x)$  qualquer

# Famílias exponenciais

uma família exponencial  $p : \mathbb{R}^d \rightarrow \mathbb{R}^+$  com parâmetros naturais  $\theta \in \mathbb{R}^d$  e dada por

$$p(x) := \exp \left( \theta^T t(x) - F(\theta) + k(x) \right); \quad x \in X$$

$t(x)$  estatística suficiente

$F(\cdot)$  função estritamente convexa

$k(x)$  qualquer

há uma expressão simples para a matriz de Fisher: 
$$I(\theta) = \frac{\partial^2 F(\theta)}{\partial \theta_i \partial \theta_j}$$

# Exemplos

Famílias exponenciais englobam muitos casos



# Exemplos

Famílias exponenciais englobam muitos casos

distribuições normais  $(x) = (x; x^2)$ ,  $(-1; 2) = (-\frac{1}{2}; \frac{1}{2})$ ,

$$F = \frac{x^2}{4} + \frac{1}{2} \log \frac{1}{2}, \quad k(x) = 0$$

# Exemplos

Famílias exponenciais englobam muitos casos

distribuição normal  $(x) = (x; x^2), (\mu; \sigma^2) = (\mu; \frac{1}{2\sigma^2})$ ,

$$F = \frac{x^2}{4\sigma^2} + \frac{1}{2} \log \frac{1}{2\sigma^2}, k(x) = 0$$

distribuição de Poisson  $(x) = x, k(x) = x!, \lambda = \log \mu, F(x) = e^{-\mu} \sum_{k=0}^x \frac{\mu^k}{k!} = e^{-\mu} e^{\mu} = e$

# Exemplos

Famílias exponenciais englobam muitos casos

distribuições normais  $(x) = (x; x^2)$ ,  $(\mu; \sigma^2) = (\mu; \frac{1}{2\sigma^2})$ ,

$$F = \frac{x^2}{4} + \frac{1}{2} \log \frac{1}{2}, \quad k(x) = 0$$

distribuições poisson  $(x) = x$ ,  $k(x) = x!$ ,  $\lambda = \log$ ,  $F(\lambda) = e^{-\lambda} = e$

gama, beta, exponencial, etc.

# Parâmetros duais

famílias exponenciais têm parametrizações duais  $r \quad F(\cdot)$

# Parâmetros duais

famílias exponenciais têm parametrizações duais  $\eta \in \mathcal{F}(\cdot)$   
e possível voltar para os parâmetros naturais via  $\eta \in \mathcal{F}(\cdot)$

# Parâmetros duais

famílias exponenciais têm parametrizações duais  $\eta \in \mathcal{F}(\eta)$   
 e possível voltar para os parâmetros naturais via  $\eta \in \mathcal{F}(\eta)$   
 onde  $\mathcal{F}(\eta) := \{h; \eta \in \mathcal{F}(\eta)\}$  e a transformada de Legendre

# Parâmetros duais

famílias exponenciais têm parametrizações duais  $\eta \in \mathcal{F}(\eta)$   
 e possível voltar para os parâmetros naturais via  $\eta \in \mathcal{F}(\eta)$   
 onde  $\mathcal{F}(\mathcal{P}) := \{h; i \in \mathcal{F}(\eta)\}$  e a transformada de Legendre  
 $\mathcal{F}(\eta) = \int_{\mathcal{X}} p(x) \eta(x) d(x)$  e a entropia de Shannon negativa

# Parâmetros duais

famílias exponenciais têm parametrizações duais  $\eta \in \mathcal{F}(\eta)$   
 e possível voltar para os parâmetros naturais via  $\eta \in \mathcal{F}(\eta)$   
 onde  $\mathcal{F}(\mathcal{P}) := \{h; i \in \mathcal{F}(\eta)\}$  e a transformada de Legendre  
 $\mathcal{F}(\eta) = \int_{\mathcal{X}} p(x) \eta(x) d(x)$  e a entropia de Shannon negativa



as parametrizações  $\eta$  e  $\theta$  de fato são duais, no sentido que:

$$\eta_i = \frac{\partial}{\partial \theta_i}; \theta^j = \frac{\partial}{\partial \eta_j} \Rightarrow g(\eta; \theta^j) = \eta_j.$$

as parametrizações  $\eta$  e  $\theta$  de fato são duais, no sentido que:

$$\eta_i = \frac{\partial \log \zeta(\theta)}{\partial \theta_i}; \quad \theta_j = \frac{\partial \log \zeta(\eta)}{\partial \eta_j} \Rightarrow \quad g(\eta; \theta) = \eta_j.$$

as parametrizações  $\eta$  e  $\theta$  de fato são duais, no sentido que:

$$\eta_i = \frac{\partial}{\partial \theta_i}; \theta_j = \frac{\partial}{\partial \eta_j} \Rightarrow g(\eta; \theta) = \eta_j \theta_j.$$

temos que  $\frac{\partial g}{\partial \eta_j}(\eta; \theta) = \theta_j$ , e  $\frac{\partial g}{\partial \theta_j}(\eta; \theta) = \eta_j$ .

as parametrizações  $\eta$  e  $\theta$  de fato são duais, no sentido que:

$$\eta_i = \frac{\partial \log \zeta(\eta)}{\partial \eta_i}; \theta_j = \frac{\partial \log \zeta(\theta)}{\partial \theta_j} \Rightarrow g(\eta; \theta_j) = \eta_j.$$

temos que  $g_{ij}(\eta) = \frac{\partial \eta_i}{\partial \theta_j}$ , e  $g^{ij}(\theta) = \frac{\partial \theta_i}{\partial \eta_j}$ .

Observação: outras famílias, como as misturas, também têm parâmetros duais

# Geodésicas duais

as parametrizações duais induzem duas geometrias  
a (e)-geometria e a (m)-geometria  
são simplesmente as geometrias planas das parametrizações

# Geodésicas duais

as parametrizações duais induzem duas geometrias  
 a (e)-geometria e a (m)-geometria  
 são simplesmente as geometrias planas das parametrizações  
 uma (e)-geodésica é uma reta nos parâmetros

# Geodésicas duais

as parametrizações duais induzem duas geometrias  
a (e)-geometria e a (m)-geometria

- são simplesmente as geometrias planas das parametrizações
- uma (e)-geodésica é uma reta nos parâmetros
- uma (m)-geodésica é uma reta nos parâmetros

# Geodésicas duais

as parametrizações duais induzem duas geometrias  
a (e)-geometria e a (m)-geometria

- são simplesmente as geometrias planas das parametrizações
- uma (e)-geodésica é uma reta nos parâmetros
- uma (m)-geodésica é uma reta nos parâmetros

as geometrias (e portanto as geodésicas) duais são descritas por uma família conexões  $\mathcal{E}_p$ ,  $\mathbb{R}^2$ , com coeficientes

$$g_{ij;k}(\theta) = E_p \left( \theta_i \theta_j + \frac{1}{2} \theta_i \theta_j \right) \theta_k^i$$



# Geodésicas duais

as parametrizações duais induzem duas geometrias localmente planas  
a (e)-geometria e a (m)-geometria

são simplesmente as geometrias planas das parametrizações  
uma (e)-geodésica é uma reta nos parâmetros  
uma (m)-geodésica é uma reta nos parâmetros

as geometrias (e portanto as geodésicas) duais são descritas por uma  
família conexões  $\Gamma$ ,  $\Gamma \in \mathbb{R}$ , com coeficientes

$$\Gamma_{ij;k}^h = E_p \left( \Gamma_{ij}^h + \frac{1}{2} \Gamma_{ij}^k \Gamma_{kl}^h \right) \Gamma_{kl}^h$$

- = 0 =) geodésicas usuais
- = 1 =) (e)-geodésicas
- = -1 =) (m)-geodésicas

# Projeções de informação

as (e)-projeções e (m)-projeções ortogonais em uma subvariedade  $S$  minimizam as entropias relativas  $D_{KL}(p|q)$  e  $D_{KL}(q|p)$

# Projeções de informação

as (e)-projeções e (m)-projeções ortogonais em uma subvariedade  $S$  minimizam as entropias relativas  $D_{KL}(p|q)$  e  $D_{KL}(q|p)$

# Distribuições discretas

os pontos são distribuições discretas  $X \in [0; 1]$

$$X = \{x_1, \dots, x_{d+1}\} \quad p$$

$$p(x_i) = p_i \in (0; 1), \quad \sum_i p_i = 1$$

# Distribuições discretas

os pontos são distribuições discretas  $X \in [0; 1]$

$$X = (x_1, \dots, x_{d+1}) \in \mathbb{P}^d$$

$$p(x_i) = p_i \in (0; 1), \quad \sum_i p_i = 1$$

essa variedade pode ser identificada com o interior do  $n$ -simplexo padrão

$$\mathbb{D}^d = \left\{ p \in \mathbb{R}^{d+1} \mid 0 < p_j < 1; \sum_i p_i = 1 \right\}$$

# Distribuições discretas

os pontos são distribuições discretas  $X \subseteq [0; 1]$

$$X = \{x_1; \dots; x_{d+1}\} \in P$$

$$p(x_i) = p_i \in (0; 1), \quad \sum_i p_i = 1$$

essa variedade pode ser identificada com o interior do  $d$ -simplexo padrão

$$d = \left\{ p \in \mathbb{R}^{d+1} \mid 0 < p_j < 1; \sum_i p_i = 1 \right\}$$

parametrização

$$\text{domínio} = \left\{ (p_1; \dots; p_d) \in \mathbb{R}_+^d \mid \sum_i p_i < 1 \right\}$$

# Distribuições discretas

os pontos são distribuições discretas  $X \subseteq [0; 1]$

$$X = \{x_1; \dots; x_{d+1}\} \in P$$

$$p(x_i) = p_i \in (0; 1), \quad \sum_i p_i = 1$$

essa variedade pode ser identificada com o interior do complexo padrão

$$D^d = \left\{ p \in \mathbb{R}^{d+1} \mid 0 < p_j < 1; \sum_{i=1}^d p_i = 1 \right\}$$

parametrização

$$\text{domínio} = \{(p_1; \dots; p_d) \in \mathbb{R}_+^d \mid \sum_{i=1}^d p_i < 1\}$$

$$(p_1; \dots; p_d) = (p_1; \dots; p_d; p_{d+1}), \quad \text{com } p_{d+1} = 1 - \sum_{i=1}^d p_i$$





matriz de Fisher:

$$g_{ij}(p) = \frac{1}{p_{d+1}} + \frac{ij}{p_i}$$

matriz de Fisher:

$$g_{ij}(p) = \frac{1}{p_{d+1}} + \frac{ij}{p_i}$$

para calcular a distância de Fisher-Rao fazemos uma reparametrização  $\alpha_i = 2^p \bar{p}_i$ ,

matriz de Fisher:

$$g_{ij}(p) = \frac{1}{p_{d+1}} + \frac{ij}{p_i}$$

para calcular a distância de Fisher-Rao fazemos uma reparametrização  $\alpha_i = 2 \sqrt{p_i}$ ,

que leva pontos  $p = (p_1, \dots, p_{d+1}) \in \mathbb{R}_+^{d+1}$  em pontos  $z$  no setor positivo da esfera

$$S_{2,+}^d = \{ z \in \mathbb{R}_+^{d+1} \mid \sum_{i=1}^{d+1} z_i^2 = 4 \}$$

matriz de Fisher:

$$g_{ij}(p) = \frac{1}{p_{d+1}} + \frac{ij}{p_i}$$

para calcular a distância de Fisher-Rao fazemos uma reparametrização  $\alpha_i = 2 \sqrt{p_i}$ ,

que leva pontos  $p = (p_1, \dots, p_{d+1}) \in \mathbb{R}_+^{d+1}$  em pontos  $z$  no setor positivo da esfera

$$S_{2,+}^d = \{z \in \mathbb{R}_+^{d+1} \mid \sum_{i=1}^{d+1} z_i^2 = 4\}$$

nessa nova parametrização, a métrica de Fisher é a métrica esférica usual de  $S_{2,+}^d \subset \mathbb{R}^{d+1}$ :  $g_{ij}(z) = \frac{\partial z}{\partial \alpha_i} \cdot \frac{\partial z}{\partial \alpha_j}$

matriz de Fisher:

$$g_{ij}(p) = \frac{1}{p_{d+1}} + \frac{ij}{p_i}$$

para calcular a distância de Fisher-Rao fazemos uma reparametrização  $\alpha_i = 2 \frac{p_i}{p_{d+1}}$ ,

que leva pontos  $p = (p_1, \dots, p_{d+1}) \in \mathbb{R}_+^{d+1}$  em pontos  $z$  no setor positivo da esfera

$$S_{2,+}^d = \{z \in \mathbb{R}_+^{d+1} \mid \sum_{i=1}^{d+1} z_i^2 = 4\}$$

nessa nova parametrização, a métrica de Fisher é a métrica esférica usual de  $S_{2,+}^d \subset \mathbb{R}^{d+1}$ :  $g_{ij}(z) = \frac{\partial z}{\partial p_i} \cdot \frac{\partial z}{\partial p_j}$

portanto a distância de Fisher-Rao entre  $p$  e  $q$  pode facilmente ser calculada como o comprimento do arco ligando  $z_p$  e  $z_q$ , que equivale a

matriz de Fisher:

$$g_{ij}(p) = \frac{1}{p_{d+1}} + \frac{ij}{p_i}$$

para calcular a distância de Fisher-Rao fazemos uma reparametrização  $z_i = 2 \sqrt{\frac{p_i}{p_{d+1}}}$ ,

que leva pontos  $p = (p_1, \dots, p_{d+1}) \in \mathbb{R}_+^{d+1}$  em pontos  $z$  no setor positivo da esfera

$$S_{2,+}^d = \{z \in \mathbb{R}_+^{d+1} \mid \sum_{i=1}^{d+1} z_i^2 = 4\}$$

nessa nova parametrização, a métrica de Fisher é a métrica esférica usual de  $S_{2,+}^d \subset \mathbb{R}^{d+1}$ :  $g_{ij}(z) = \frac{\partial z_i}{\partial p_i} \frac{\partial z_j}{\partial p_j}$

portanto a distância de Fisher-Rao entre  $p$  e  $q$  pode facilmente ser calculada como o comprimento do arco ligando  $z_p$  e  $z_q$ , que equivale a

$$d_{FR}(p; q) = 2 \arccos \sqrt{\sum_{i=1}^{d+1} p_i q_i}$$

um fato interessante é que o comprimento da corda ligando  $z_p$  e  $z_q$  fornece uma boa aproximação:

um fato interessante é que o comprimento da corda ligando  $z_p$  e  $z_q$  fornece uma boa aproximação:

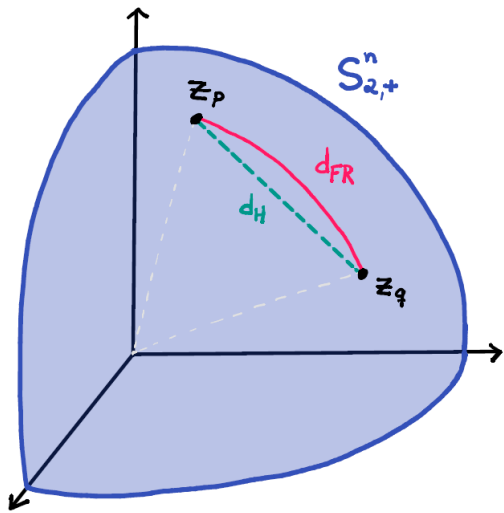
$$z_p - z_q \approx 2 \sqrt{\sum_{i=1}^{n-1} p_i \bar{p}_i - p \bar{q}^2}$$



- um fato interessante é que o *comprimento da corda* ligando  $Z_p$  e  $Z_q$  fornece uma boa aproximação:

$$Z_p \quad Z_q = 2 \sqrt{\sum_{i=1}^{d+1} \left( \sqrt{\frac{p_i}{q_i}} - \sqrt{\frac{q_i}{p_i}} \right)^2}$$

- essa distância, sem o fator 2, é chamada *distância de Hellinger*  $d_H$



# Parte III

## Aplicações

# Funções perda para classificadores

problemas de classificação aprendizado de máquina:

# Funções perda para classificadores

problemas de classificação aprendizado de máquina:  
temos uma família parametrizada de funções  $f : X \rightarrow \mathbb{R}^K$

# Funções perda para classificadores

problemas de classificação aprendizado de máquina:  
 temos uma família parametrizada de funções  $f : X \rightarrow \mathbb{R}^K$   
 $X \subset \mathbb{R}^n$  e o espaço dos vetores de características (ex: imagens)

# Funções perda para classificadores

problemas de classificação aprendizado de máquina:  
 temos uma família parametrizada de funções  $f : X \rightarrow \mathbb{R}^K$   
 $X \subset \mathbb{R}^n$  e o espaço dos vetores de características (ex: imagens)  
 $K$  e o número de classes (ex: cachorro, gato, etc.)

# Funções perda para classificadores

problemas de classificação aprendizado de máquina:

temos uma família parametrizada de funções  $f : X \rightarrow \mathbb{R}^K$

$X \subset \mathbb{R}^n$  e o espaço dos vetores de características (ex: imagens)

$K$  e o número de classes (ex: cachorro, gato, etc.)

$z = f(x)$  e chamado vetorescore



# Funções perda para classificadores

problemas de classificação aprendizado de máquina:

temos uma família parametrizada de funções  $f : X \rightarrow \mathbb{R}^K$

$X \subset \mathbb{R}^n$  e o espaço dos vetores de características (ex: imagens)

$K$  e o número de classes (ex: cachorro, gato, etc.)

$z = f(x)$  e chamado vetorescore

$\theta$  são os parâmetros da máquina (geralmente dados por uma rede neural)

# Funções perda para classificadores

problemas de classificação aprendizado de máquina:  $n$  0

temos uma família parametrizada de funções  $f : X \rightarrow \mathbb{R}^K$   $2$

$X \subset \mathbb{R}^n$  e o espaço dos vetores de características (ex: imagens)

$K$  e o número de classes (ex: cachorro, gato, etc.)

$z = f(x)$  e chamado vetorescore

$2$  são os parâmetros da máquina (geralmente dados por uma rede neural)

transformamos o vetor escore em um vetor de probabilidades através da função softmax  $\sigma : \mathbb{R}^K \rightarrow \mathbb{R}^{K-1}$  dada em coordenadas por

$$(z)_i = \frac{e^{z_i}}{\sum_{i=1}^K e^{z_i}}$$

# Funções perda para classificadores

problemas de classificação aprendizado de máquina:  $n$  0

temos uma família parametrizada de funções  $f : X \rightarrow \mathbb{R}^K$   $2$

$X \subset \mathbb{R}^n$  e o espaço dos vetores de características (ex: imagens)

$K$  e o número de classes (ex: cachorro, gato, etc.)

$z = f(x)$  e chamado vetorescore

$2$  são os parâmetros da máquina (geralmente dados por uma rede neural)

transformamos o vetor escore em um vetor de probabilidades através da função softmax  $\sigma : \mathbb{R}^K \rightarrow \mathbb{R}^{K \times 1}$  dada em coordenadas por

$$(z)_i = \frac{e^{z_i}}{\sum_{i=1}^K e^{z_i}}$$

podemos interpretar  $(z)_i$  como a probabilidade do vetor pertencer a classe  $i$

treinamento supervisionado: somos fornecidos com um conjunto de  
 treinamento  $(x_i; y_i)_{i=1}^m \times \{1, \dots, K\}$

treinamento supervisionado: somos fornecidos com um conjunto de  
 treinamento  $(x_i; y_i)_{i=1}^m \times \{1, \dots, K\}$   
 tomamos uma função perda  $\ell : \mathbb{R}^K \rightarrow \mathbb{R}_+$

treinamento supervisionado: somos fornecidos com um conjunto de treinamento  $(x_i; y_i)_{i=1}^m \subset \mathbb{R}^K \times \mathbb{R}^K$

tomamos uma função perda  $L: \mathbb{R}^K \times \mathbb{R}^K \rightarrow \mathbb{R}_+$

o problema de aprendizado de máquina consiste em minimizar a perda média do conjunto de treinamento:

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^m L(\theta; x_i, y_i)$$

treinamento supervisionado: somos fornecidos com um conjunto de treinamento  $(x_i; y_i)_{i=1}^m$   $X = \{x_1, \dots, x_m\}$   $f: \mathbb{R}^d \rightarrow \mathbb{R}$

tomamos uma função perda  $L: \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}_+$

o problema de aprendizado de máquina consiste em minimizar a perda média do conjunto de treinamento:

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^m L(f(x_i); y_i)$$

isso costuma ser feito através do método do gradiente

funções perda mais usadas:  
entropia cruzada:  $h(p; q) = - \sum_i p_i \log \frac{1}{q_i}$



funções perda mais usadas:

entropia cruzada:  $h(p; q) = \sum_i p_i \log \frac{1}{q_i}$

perda quadrática:  $L(p; q) = \|p - q\|_2^2$

funções perda mais usadas:

$$\text{entropia cruzada: } h(p; q) = \sum_i p_i \log \frac{1}{q_i}$$

$$\text{perda quadrática: } L(p; q) = k \|p - q\|_2^2$$

nossa proposta: usar as perdas geométrico-informacionais no simplexo, dadas pelo quadrado das distâncias apresentadas:

funções perda mais usadas:

$$\text{entropia cruzada: } h(p; q) = \sum_i p_i \log \frac{1}{q_i}$$

$$\text{perda quadrática: } L(p; q) = k \|p - q\|_2^2$$

nossa proposta: usar as perdas geométrico-informacionais no simplexo, dadas pelo quadrado das distâncias apresentadas:

$$4L_{\text{SFR}} = d_{\text{FR}}^2(p; q) = 4 \arccos \left( \sum_{i=1}^K p_i \sqrt{p_i q_i} \right)^2$$

$$L_{\text{SH}} = d_{\text{H}}^2(p; q) = \sum_{i=1}^K p_i \left| \frac{p_i}{q_i} - 1 \right|^2$$

funções perda mais usadas:

$$\text{entropia cruzada: } h(p; q) = - \sum_i p_i \log \frac{1}{q_i}$$

$$\text{perda quadrática: } L(p; q) = k \sum_i (p_i - q_i)^2$$

nossa proposta: usar as perdas geométrico-informacionais no simplexo, dadas pelo quadrado das distâncias apresentadas:

$$4L_{\text{SFR}} = d_{\text{FR}}^2(p; q) = 4 \arccos \left( \sum_{i=1}^K p_i \sqrt{\frac{p_i}{q_i}} \right)^2$$

$$L_{\text{SH}} = d_{\text{H}}^2(p; q) = \sum_{i=1}^K \left( \sqrt{\frac{p_i}{q_i}} - \sqrt{\frac{q_i}{p_i}} \right)^2$$

este é um trabalho em conjunto com H.K. Miyamoto e S.I.R. Costa, submetido para o ISIT 2022 (International Symposium on Information Theory)<sup>2</sup>

<sup>2</sup>Henrique K. Miyamoto, Fabio C. C. Meneghetti e Sueli I. R. Costa.

observamos que existem relações assintóticas e desigualdades entre perdas que introduzimos e a perda da entropia cruzada

observamos que existem relações assintóticas e desigualdades entre perdas que introduzimos e a perda da entropia cruzada



# Resultados

(a) MNIST

(b) Banco de dados sintético

Figura: Acurácia dos aprendizados com diferentes funções perda.

# Resultados com rudo

(alguns rotulos do conjunto de treinamento recebem a classe errada)

(a) MNIST, = 0:3

(b) MNIST, = 0:5



(a) Sintético,  $\rho = 0:3$

(b) Sintético,  $\rho = 0:5$

possivelmente uma das razões das perdas de Fisher-Rao e Helling  
terem boa performance no caso ruidoso seja por elas serem limitadas

possivelmente uma das razões das perdas de Fisher-Rao e Helling  
terem boa performance no caso ruidoso seja por elas serem limitadas  
este é um trabalho em andamento

# Distribuições de probabilidade enroladas no toro

um reticulado posto-completo e um subconjunto de  $\mathbb{R}^n$  gerado por combinações lineares inteiras de uma base  $\{b_1, \dots, b_n\}$

# Distribuições de probabilidade enroladas no toro

um reticulado posto-completo e um subconjunto de  $\mathbb{R}^n$  gerado por combinações lineares inteiras de uma base  $b_1, \dots, b_n$  de  $\mathbb{R}^n$ . Assim, o toro enrolado por  $\Gamma$  como

$$T := \mathbb{R}^n / \Gamma = \mathbb{R}^n / \left[ \sum_{i=1}^n x_i b_i \mid x_i \in \mathbb{Z} \right]$$

# Distribuições de probabilidade enroladas no toro

um reticulado posto-completo e um subconjunto de  $\mathbb{R}^n$  gerado por combinações lineares inteiras de uma base  $\{b_1, \dots, b_n\}$  de  $\mathbb{R}^n$ . Assim, o toro enrolado por  $\Gamma$  como

$$T := \mathbb{R}^n / \Gamma = \{ [x] = x + \Gamma \mid x \in \mathbb{R}^n \}$$

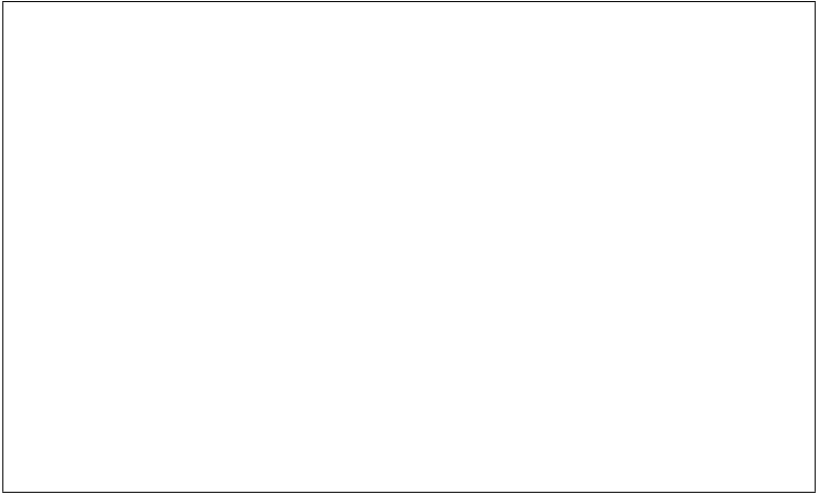
temos uma projecão canónica  $\pi : \mathbb{R}^n \rightarrow T$ ,  $\pi(x) = [x]$

# Distribuições de probabilidade enroladas no toro

um reticulado posto-completo e um subconjunto de  $\mathbb{R}^n$  gerado por combinações lineares inteiras de uma base  $b_1, \dots, b_n$  de  $\mathbb{R}^n$ . Assim, o toro enrolado por  $\Gamma$  como

$$T := \mathbb{R}^n / \Gamma = \mathbb{R}^n / \langle [x] = x + \sum_{i=1}^n x_i b_i \mid x_i \in \mathbb{Z} \rangle$$

temos uma projeção canônica  $\pi : \mathbb{R}^n \rightarrow T$ ,  $\pi(x) = [x]$ .  
 dada uma distribuição de probabilidade  $P$  em  $\mathbb{R}^n$ , podemos enrolá-la no toro com um empurro via  $\pi$ , isto é,  $P := (\pi)_* P$





se uma distribuição  $P$  em  $\mathbb{R}^n$  tem densidade  $p(x)$ ,  $x \in \mathbb{R}^n$ , então a distribuição enrolada tem densidade

$$p(x) = \sum_{k \in \mathbb{Z}^n} p(x + k)$$

sobre  $T$  ou uma região fundamental (ex: região de Voronoi)

se uma distribuição  $P$  em  $\mathbb{R}^n$  tem densidade  $p(x)$ ,  $x \in \mathbb{R}^n$ , então a distribuição enrolada tem densidade

$$p(x) = \frac{1}{2} (p(x + \cdot) + p(x - \cdot))$$

sobre  $T$  ou uma região fundamental (ex: região de Voronoi) assim, a partir de um modelo estatístico  $p$  em  $\mathbb{R}^n$  obtemos um modelo estatístico  $p$  em  $T$ .

se uma distribuição  $p$  em  $\mathbb{R}^n$  tem densidade  $p(x)$ ,  $x \in \mathbb{R}^n$ , então a distribuição enrolada tem densidade

$$p(x) = \sum_{k \in \mathbb{Z}^n} p(x + k)$$

sobre  $T$  ou uma região fundamental (ex: região de Voronoi) assim, a partir de um modelo estatístico  $p$  em  $\mathbb{R}^n$  obtemos um modelo estatístico  $p; \kappa$  em  $T$ .

ex: gaussianas multivariadas

$$p; \kappa; (x) = \frac{1}{(2\pi)^d |\det K|} e^{-\frac{1}{2}(x - \mu)^T K^{-1}(x - \mu)}$$

uma propriedade central dessas distribuições e que para variância crescente elas se aproximam da distribuição uniforme  $\mu(\alpha) = \frac{1}{j \det j}$

uma propriedade central dessas distribuições e que para variância crescente elas se aproximam da distribuição uniforme  $\mu(\alpha) = \frac{1}{j \det j}$



gostaramos de estudar a geometria dessas distribuições



gostaramos de estudar a geometria dessas distribuições  
em termos da geometria de Fisher-Rao





gostaramos de estudar a geometria dessas distribuições  
em termos da geometria de Fisher-Rao  
e em termos de medidas de divergência (Kulback-Leibler,  
f-divergências, normal  $\mathbb{S}^p$ , etc.)

gostaramos de estudar a geometria dessas distribuçoes  
 em termos da geometria de Fisher-Rao  
 e em termos de medidas de divergência (Kulback-Leibler,  
 f-divergências, normal<sup>p</sup>, etc.)

já existe alguma pesquisa sobre distribuçoes desse tipo em termos  
 geometria de Wasserstein<sup>3</sup>

---

<sup>3</sup>Anton Mallasto e Aasa Feragen. "Optimal Transport Distance between Wrapped Gaussian Distributions". Em: 38th MaxEnt. 2018

gostaramos de estudar a geometria dessas distribuições  
 em termos da geometria de Fisher-Rao  
 e em termos de medidas de divergência (Kulback-Leibler,  
 f-divergências, normal<sup>P</sup>, etc.)

já existe alguma pesquisa sobre distribuições desse tipo em termos  
 geometria de Wasserstein<sup>3</sup>

nossa motivação: o fator de achatamento<sup>α</sup> (ness factor) e a  
 distância<sup>L<sup>1</sup></sup> entre uma distribuição<sup>P</sup> e a uniforme

---

<sup>3</sup>Anton Mallasto e Aasa Feragen. "Optimal Transport Distance between Wrapped Gaussian Distributions". Em: 38th MaxEnt. 2018

gostaramos de estudar a geometria dessas distribuições em termos da geometria de Fisher-Rao e em termos de medidas de divergência (Kulback-Leibler, f-divergências, normal<sup>P</sup>, etc.)

já existe alguma pesquisa sobre distribuições desse tipo em termos geometria de Wasserstein<sup>3</sup>

nossa motivação: o fator de achatamento (ness factor) e a distância  $L^1$  entre uma distribuição<sup>P</sup> e a uniforme

ele é um parâmetro importante para construir códigos que atingem capacidade no canal AWGN e para garantir segredo no canal Wiretap

---

<sup>3</sup>Anton Mallasto e Aasa Feragen. "Optimal Transport Distance between Wrapped Gaussian Distributions". Em: 38th MaxEnt. 2018

gostaramos de estudar a geometria dessas distribuições em termos da geometria de Fisher-Rao e em termos de medidas de divergência (Kulback-Leibler, f-divergências, normal<sup>P</sup>, etc.)

já existe alguma pesquisa sobre distribuições desse tipo em termos geometria de Wasserstein<sup>3</sup>

nossa motivação: o fator de achatamento (ness factor) e a distância  $L^1$  entre uma distribuição<sup>P</sup> e a uniforme

é um parâmetro importante para construir códigos que atingem capacidade no canal AWGN e para garantir segredo no canal Wiretap. Queremos entender se o fator de achatamento medido com outras divergências também tem comportamento interessante

---

<sup>3</sup>Anton Mallasto e Aasa Feragen. "Optimal Transport Distance between Wrapped Gaussian Distributions". Em: 38th MaxEnt. 2018

gostamos de estudar a geometria dessas distribuições em termos da geometria de Fisher-Rao e em termos de medidas de divergência (Kulback-Leibler, f-divergências, normal<sup>P</sup>, etc.)

já existe alguma pesquisa sobre distribuições desse tipo em termos de geometria de Wasserstein<sup>3</sup>

nossa motivação: o fator de achatamento (ness factor) e a distância  $L^1$  entre uma distribuição<sup>P</sup> e a uniforme

é um parâmetro importante para construir códigos que atingem capacidade no canal AWGN e para garantir segredo no canal Wiretap. Queremos entender se o fator de achatamento medido com outras divergências também tem comportamento interessante. Este tema está diretamente conectado ao mestrado do aluno

<sup>3</sup>Anton Mallasto e Aasa Feragen. "Optimal Transport Distance between Wrapped Gaussian Distributions". Em: 38th MaxEnt. 2018

<sup>4</sup>Fabio C. C. Meneghetti. "Reticulados: um estudo de alguns parâmetros relevantes para aplicações em criptografia". 2020

# Parte IV

## Futuro

# Futuro

queremos continuar estudando alguns aspectos teóricos



# Futuro

queremos continuar estudando alguns aspectos teóricos  
 a relação entre a estrutura dualmente plana de  $(M, g, \omega, \mu)$  e  
 as geometrias simplética e Kähler

# Futuro

queremos continuar estudando alguns aspectos teóricos  
 a relação entre a estrutura dualmente plana de  $(M, g; r; r)$  e  
 as geometrias simplética e Kähler  
 extensões da geometria da informação para espaços de dimensão  
 infinita (ex: estrutura de Pistone-Sempi)

# Futuro

queremos continuar estudando alguns aspectos teóricos

a relação entre a estrutura dualmente plana de  $(M, g; r; r)$  e

as geometrias simplética e Kähler

extensões da geometria da informação para espaços de dimensão infinita (ex: estrutura de Pistone-Sempi)

entender se há relação entre nossa proposta de funções perda, e as

-Divergências, e também com o método do gradiente natural

# Futuro

queremos continuar estudando alguns aspectos teóricos

a relação entre a estrutura dualmente plana de  $(M, g; r; r)$  e

as geometrias simplética e Kähler

extensões da geometria da informação para espaços de dimensão infinita (ex: estrutura de Pistone-Sempi)

entender se há relação entre nossa proposta de funções perda, e as

-Divergências, e também com o método do gradiente natural

formalizar a teoria das distribuições gaussianas enroladas no toro

# Futuro

queremos continuar estudando alguns aspectos teóricos

a relação entre a estrutura dualmente plana de  $(M, g; r; r^{-1})$  e

as geometrias simplética e Kähler

extensões da geometria da informação para espaços de dimensão infinita (ex: estrutura de Pistone-Sempi)

entender se há relação entre nossa proposta de funções perda, e as

-Divergências, e também com o método do gradiente natural

formalizar a teoria das distribuições gaussianas enroladas no toro

procurar relações entre reticulados diferentes (ex: se  $B = \mathbb{Z}^n$ , então

$$p; \kappa; (x) = \frac{1}{\det} p_{\tilde{\kappa}; \mathbb{Z}^n}(B^{-1}x), \quad \tilde{\kappa} = B^{-1}, \quad \tilde{\kappa} = B^{-1}KB^{-t}$$

# Futuro

queremos continuar estudando alguns aspectos teóricos

a relação entre a estrutura dualmente plana de  $\mathbb{A}^m(\mathbb{C}; g; r; r')$  e

as geometrias simplética e Kähler

extensões da geometria da informação para espaços de dimensão infinita (ex: estrutura de Pistone-Sempi)

entender se há relação entre nossa proposta de funções perda, e as

-Divergências, e também com o método do gradiente natural

formalizar a teoria das distribuições gaussianas enroladas no toro

procurar relações entre reticulados diferentes (ex: se  $B = \mathbb{Z}^n$ , então

$p; \kappa; (x) = \frac{1}{\det B} p_{\tilde{\kappa}; \mathbb{Z}^n}(B^{-1}x)$ ,  $\tilde{\kappa} = B^{-1} \kappa$ ,  $\tilde{\kappa} = B^{-1}KB^{-t}$ )

reticulados duais parecem ter relação com a transformada de Fourier distribuição

# Livros

- [1] Shun'ichi Amari e Hiroshi Nagaoka **Methods of information geometry** Trad. por Daishi Harada. Translations of mathematical monographs. American Mathematical Society, 2007.
- [3] Nihat Ay, Jürgen Jost, Hồng Vân Lê e Lorenz Schwachhofer. **Information Geometry** Springer International Publishing, 2017.
- [4] Ovidiu Calin e Constantin Udriste **Geometric Modeling in Probability and Statistics** Springer International Publishing, 2014.

# Artigos

- [2] Colin Atkinson e Ann F. S. Mitchell. "Rao's Distance Measure". Em: Sankhya: The Indian Journal of Statistics, Series(A)981).
- [6] Ahmet Demirkaya, Jiasi Chen e Samet Oymak. "Exploring the Role of Loss Functions in Multiclass Classification". Em: 54th CISS. 2020
- [7] Anton Mallasto e Aasa Feragen. "Optimal Transport Distance between Wrapped Gaussian Distributions". Em: 38th MaxEnt. 2018
- [8] Fabio C. C. Meneghetti. "Reticulados: um estudo de alguns parâmetros relevantes para aplicações em criptografia". 2020.
- [10] Frank Nielsen. "An Elementary Introduction to Information Geometry". Em:Entropy (2020).
- [12] Julianna Pinele, Joao E. Strapasson e Sueli I. R. Costa. "The Fisher-Rao Distance between Multivariate Normal Distributions: Special Cases, Bounds and Applications". Em:Entropy (2020).



- [5] Alberto Cena e Giovanni Pistone. “Exponential statistical manifold”. Em: *Annals of the Institute of Statistical Mathematics* (2006).
- [11] Tomonori Noda. “Symplectic Structures on Statistical Manifolds”. Em: *J. Aust. Math. Soc.* (2011).
- [13] Rui F. Vigelis, Luiza H. F. De Andrade e Charles C. Cavalcante. “Properties of a Generalized Divergence Related to Tsallis Generalized Divergence”. Em: *IEEE Transactions on Information Theory* (2020).