

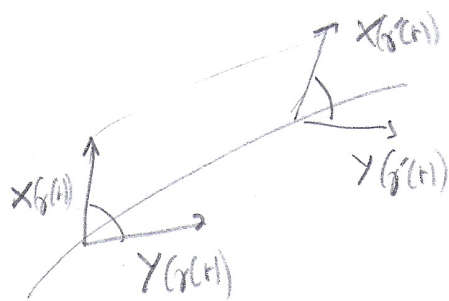
# Riemannian connection

There are many possible choice for affine connections  
 If one has a notion of metric, a natural choice is  
 to choose the Riemannian connection.

\* symmetry:  $\Gamma_{ij}^k = \Gamma_{ji}^k$  (torsion free)

\* preserve inner product by parallel transport:

$$? \quad \frac{d}{dt} \langle X(\gamma(t)), Y(\gamma(t)) \rangle = \langle \underset{\substack{\uparrow \\ \text{covariant derivative}}}{\nabla_{\dot{\gamma}(t)} X}, Y \rangle + \langle X, \underset{\substack{\uparrow \\ \text{covariant derivative}}}{\nabla_{\dot{\gamma}(t)} Y} \rangle$$



parallel transport:

$$\nabla_{\dot{\gamma}} X = \underbrace{\left\{ \dot{X}^k + \dot{\gamma}^i X^k \Gamma_{ij}^k \right\}}_0 \partial_k|_{\gamma(t)}$$

$X, Y$  parallel translated  $\Rightarrow \frac{d}{dt} \langle X(\gamma(t)), Y(\gamma(t)) \rangle = 0$

Requirement to be true for any path and  $X, Y, \dot{\gamma}(t) = Z$

$$\partial_k \langle \partial_i, \partial_j \rangle = \langle \nabla_{\partial_k} \partial_i, \partial_j \rangle + \langle \partial_i, \nabla_{\partial_k} \partial_j \rangle$$

$$\Rightarrow \text{local: } \begin{cases} \partial_k g_{ij} = \Gamma_{ki,j}^k + \Gamma_{kj,i}^k \\ \Gamma_{ij,k}^k = \frac{1}{2} (\partial_i g_{jk} + \partial_j g_{ki} - \partial_k g_{ij}) \end{cases}$$

## Information - based connection

Natural notion Kullback Leibler divergence

$$D[p \parallel q] = \int dx p(x) \log \frac{p(x)}{q(x)}$$

Problem: not symmetric, do not satisfy triangular inequality

However:  $D(p \parallel q) \geq 0$       $D(p \parallel q) (=) p = q$   
                                   $\uparrow \uparrow$   
                                  smooth

Notations      $D(\partial_i | p \parallel p') = \partial_i D(p \parallel p')$

$$D(\partial_i | p \parallel \partial_j | p') = \partial_i \partial_j (p \parallel p') \dots$$

Diagonal constant  $\Rightarrow D[\partial_i \parallel] = D[ \parallel \partial_j ] = 0$

$$\Rightarrow D[\partial_i \partial_j \parallel] + D[\partial_i \parallel \partial_j] = 0$$

$$\Rightarrow \underbrace{D[\partial_i \partial_j \parallel]}_{\underbrace{\quad}_{g_{ij}(p)}} = - \underbrace{D[\partial_i \parallel \partial_j]}_{\underbrace{\quad}_{g_{ij}(p)}} = D[\parallel \partial_i \partial_j]_{\underbrace{\quad}_{g_{ij}(p)}}$$

$$D(p \parallel q) = \frac{1}{2} g_{ij}(q) \Delta z^i \Delta z^j + o(\|\Delta z\|^2) \quad \Delta z^i = z^i(p) - z^i(q)$$

$\uparrow$   
metric!

Connection  $\rightarrow$  3rd order      $\Gamma_{ij,k} = -D[\partial_i \partial_j \parallel \partial_k]$

$$D[p||q] = \frac{1}{2} g_{ij}(q) \Delta z^i \Delta z^j + \frac{1}{6} h_{ijk}(q) \Delta z^i \Delta z^j \Delta z^k + o(\|\Delta z\|^3)$$

$$h_{ijk} = D[\partial_i \partial_j \partial_k ||] = \partial_i [D[\partial_j \partial_k ||]] - D[\partial_j \partial_k || \partial_i] \\ = \partial_i g_{jk} + \Gamma_{jk,i}$$

$D^*[q||p] = D[p||q] \rightarrow D^*$  natural dual function

$$= \frac{1}{2} g_{ij}(p) \Delta z^i \Delta z^j + \frac{1}{6} h_{ijk}^*(p) \Delta z^i \Delta z^j \Delta z^k + o(\|\Delta z\|^3) \\ \frac{1}{2} \left( g_{ij}(q) \Delta z^i \Delta z^j + \partial_k g_{ij}(p) \Delta z^i \Delta z^j \Delta z^k \right)$$

$$h_{ijk}^* = D[|| \partial_i \partial_j \partial_k] = \partial_i g_{jk} + \Gamma_{jk,i}^*, \quad \Gamma_{jk,i}^* = -D[\partial_i || \partial_j \partial_k]$$

$\Rightarrow$  Dual relation:  $\partial_k g_{ij} = \Gamma_{ki,j} + \Gamma_{jk,i}^*$

$$\partial_k g_{ij} = -\partial_k D[\partial_i || \partial_j] = -D[\partial_k \partial_i || \partial_j] - D[\partial_i || \partial_k \partial_j]$$

Almost Riemannian:  $\Gamma_{ki,j} = \Gamma_{ki,j}^*$  self-dual

F-divergence

convex, smooth

$$D_F(p||q) = 0 \Leftrightarrow p = q \quad \textcircled{4}$$

if F strictly convex at 1

$$D_F(p||q) = \int p(x) F\left(\frac{q(x)}{p(x)}\right) dx$$

↑

$$\text{Jensen inequality: } D_F(p||q) \geq F\left(\int p(x) \frac{q(x)}{p(x)} dx\right) = F(1) \stackrel{\text{def}}{=} 0$$

$$\text{Monotonicity: } \kappa(y|x) \rightarrow p_\kappa(y) = \int \kappa(y|x) p(x) dx$$

↓

Convexity

$$q_\kappa(y) = \int \kappa(y|x) q(x) dx$$

$$D_F(p||q) = \int dx \int dy \kappa(y|x) p(x) F\left(\frac{q(x)}{p(x)}\right) dx$$

$$p_\kappa(x|y) = \frac{\kappa(y|x)p(x)}{p_\kappa(y)} = \int dy p_\kappa(y) \int dx p_\kappa(x|y) F\left(\frac{q(x)}{p(x)}\right)$$

$$\text{Jensen} \rightarrow \geq \int dy p_\kappa(y) F\left(\int dx p_\kappa(x|y) \frac{q(x)}{p(x)}\right)$$

$$= \int dy p_\kappa(y) F\left(\frac{q_\kappa(y)}{p_\kappa(y)}\right) = D_F(p_\kappa||q_\kappa)$$

$$D_F(p||q) = D_F(p_\kappa||q_\kappa) \Leftrightarrow \kappa \text{ sufficient statistics } p, q$$

$$\Leftrightarrow p_\kappa(x|y) = q_\kappa(x|y)$$

$$\kappa(y|x) = \int_{F(x)} \delta(y)$$

(strictly convex)

$$F^\alpha(u) = \begin{cases} \frac{4}{1-\alpha^2} \left\{ 1 - u^{(1+\alpha)/2} \right\} & \alpha \neq \pm 1 \\ u \log u & \\ -\log u & \end{cases} \text{KL}$$

$\alpha$ -divergence

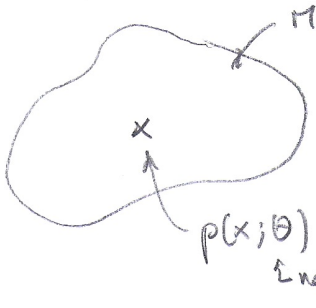
$$D^{(0)}(p||q) = \int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx$$

$$\sqrt{D^{(0)}(p||q)}$$

Hellinger distance

Fisher metric

$\theta^i$ : coordinates.



$$g_{ij}(\theta) = \mathbb{E}_\theta \left[ \partial_i \log p \partial_j \log p \right] \leftarrow \text{covariance 2.}$$

$$= \int dx \frac{\partial_i p}{\sqrt{p}} \frac{\partial_j p}{\sqrt{p}}$$

Statistical model:  $p(x, \theta) = \begin{cases} \theta_i, & x=i \\ \theta_0 = 1 - \sum_{i=1}^m \theta_i \end{cases}$

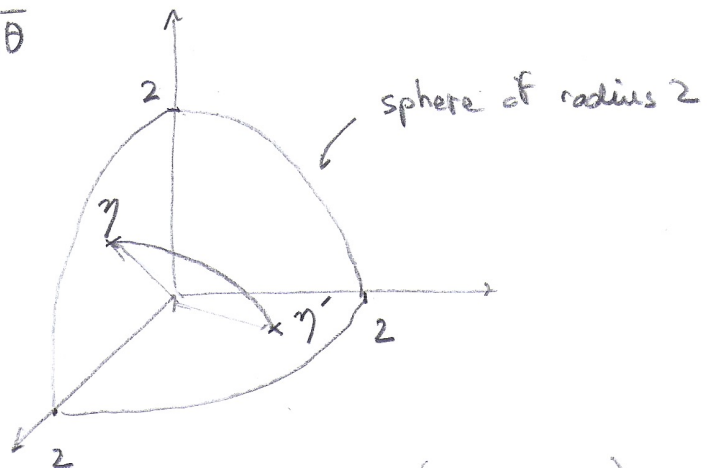
$$g_{ij}(\theta) = \sum_k \theta_k \partial_i \log \theta_k \partial_j \log \theta_k + \theta_0 \partial_i \log \theta_0 \partial_j \log \theta_0$$

$$= \sum_k \theta_k \frac{\delta_{ik}}{\theta_k} \frac{\delta_{jk}}{\theta_k} + \theta_0 \frac{(-1)}{\theta_0} \frac{(-1)}{\theta_0} = \frac{\delta_{ij}}{\theta_i} + \frac{1}{1 - \sum \theta_i} \leftarrow m \text{ variables}$$

$g_{ij}$  can be seen as an induced metric:

$$g_{ij}(\theta) = \sum_{k=1}^n \frac{\partial_i \theta_k}{\sqrt{\theta_k}} \frac{\partial_j \theta_k}{\sqrt{\theta_k}} + \frac{\partial_i \theta_0}{\sqrt{\theta_0}} \frac{\partial_j \theta_0}{\sqrt{\theta_0}} = \langle 2 \partial_i \sqrt{\theta}, 2 \partial_j \sqrt{\theta} \rangle_{n+1}$$

$$\sqrt{\theta} = \begin{pmatrix} \sqrt{\theta_0} \\ \sqrt{\theta_1} \\ \vdots \\ \sqrt{\theta_n} \end{pmatrix}, \quad \eta = 2\sqrt{\theta}$$



$$\langle \eta, \eta' \rangle_{n+1} = 2 \sum_{k=0}^n \sqrt{\theta_k} \sqrt{\theta'_k} \Rightarrow L(\theta, \theta') = 2 \cos^{-1} \left( \sum_{k=0}^{n+1} \sqrt{\theta_k \theta'_k} \right)$$

Hellinger distance:  $\sum_k (\sqrt{\theta_k} - \sqrt{\theta'_k})^2 = 2 \left( 1 - \sum_{k=0} \sqrt{\theta_k \theta'_k} \right) = 2 \sin^2 (L(\theta, \theta') / 4)$

# Property of Fisher metric

Sufficient statistics : model  $p(x|\theta)$

$$F: X \rightarrow Y \quad \left| \quad r(x|\theta) = \frac{p(x|\theta)}{q(F(x)|\theta)} \rightarrow \text{does not depend on } \theta \text{ for all } x$$

$$x \mapsto F(x) = y$$

$$p(x|y, \theta) = r(x|\theta) \int_{F(x)}(y)$$

$$k(y|x) \rightarrow g_k, \quad g(\theta) - g_k(\theta) \leftarrow \text{positive semidefinite}$$

$$\Delta g(\theta)$$

$$\Delta g(\theta) = \mathbb{E}_\theta [\partial_i \log r(x|\theta) \partial_j \log r(x|\theta)] \text{ if deterministic}$$

$$= \mathbb{E}_\theta [\text{Cov}[\partial_i \log p(x|\theta) \partial_j \log p(x|\theta)] | y]$$

## Monotonicity / Chain rule

Additivity:  $g(\theta) = g_1(\theta) + g_2(\theta)$  if  $p(x|\theta) = p(x_1|\theta)p(x_2|\theta)$

Convexity:  $g_\lambda(\theta) \leq \lambda g_1(\theta) + (1-\lambda)g_2(\theta)$

## Cramer Rao inequality

unbiased estimator  $\mathbb{E}_\theta [\hat{\theta}(x)] = \theta \quad \forall \theta$

$$V_\theta = \mathbb{E}_\theta [(\hat{\theta}^i(x) - \theta^i)(\hat{\theta}^j(x) - \theta^j)]$$

$$V_\theta - \hat{g}^{-1}(\theta) \geq 0 \quad \text{positive semidefinite}$$

## $\alpha$ -connection

(7)

$$\Gamma_{ij,k}^{(\alpha)} = \mathbb{E}_{\theta} \left[ \left( \partial_i \partial_j \log p_{\theta} + \frac{1-\alpha}{2} \partial_i \log p_{\theta} \partial_j \log p_{\theta} \right) \partial_k \log p_{\theta} \right]$$

$$\langle \nabla_{\partial_i}^{(\alpha)} \partial_j, \partial_k \rangle = \Gamma_{ij,k}^{(\alpha)}$$

$$\nabla^{(\alpha)} = (1-\alpha) \nabla^{(0)} + \alpha \nabla^{(1)} + \frac{1+\alpha}{2} \nabla^{(1)} + \frac{1-\alpha}{2} \nabla^{(-1)}$$

$$\partial_k g_{ij} = \Gamma_{ki,j}^{(0)} + \Gamma_{kj,i}^{(0)} \quad \rightarrow \text{Riemannian metric for the Hellinger distance.}$$

## Flat statistical model

Exponential family:  $p(x|\theta) = \exp(C(x) + \sum \theta^i F_i(x) - \Psi(\theta))$   
 $\Psi(\theta) = \log \int \exp[C(x) + \sum \theta^i F_i(x)] dx$

Mixture family:  $p(x|\theta) = C(x) + \sum \theta^i F_i(x)$

Exp:  $C(x) = 0$      $F_i(x) = \begin{cases} 1 & x=i \\ 0 & \end{cases}$  ,     $\theta^i = \log \frac{p_i}{p_0}$   
 $\Psi(\theta) = \log \left( 1 + \sum_i e^{\theta^i} \right)$

$$\left. \begin{aligned} \partial_i \log p(x|\theta) &= F_i(x) - \partial_i \Psi(\theta) \\ \partial_i \partial_j \log p(x|\theta) &= -\partial_i \partial_j \Psi(\theta) \end{aligned} \right\} \begin{array}{l} \text{Exponential family} \\ \text{are 1-flat} \end{array}$$

$\theta^i$  affine coordinate     $\nabla^{(1)}$  exponential connection

## Property of Fisher metric

Sufficient statistics : model  $p(x|\theta)$

$$F: X \rightarrow Y \quad \left| \quad r(x|\theta) = \frac{p(x|\theta)}{q(F(x)|\theta)} \rightarrow \begin{array}{l} \text{does not depend} \\ \text{on } \theta \text{ for all } x \end{array}$$

$$x \mapsto F(x) = y$$

$$p(x|y, \theta) = r(x|\theta) \int_{F(x)} q(y)$$

$$k(y|x) \rightarrow g_k, \quad \begin{array}{l} g(\theta) - g_k(\theta) \leftarrow \text{positive semi-definite} \\ \text{"} \\ \Delta g(\theta) \end{array}$$

$$\Delta g(\theta) = \mathbb{E}_\theta [\partial_i \log r(x|\theta) \partial_j \log r(x|\theta)] \text{ if deterministic}$$

$$= \mathbb{E}_\theta [\text{Cov}[\partial_i \log p(x|\theta) \partial_j \log p(x|\theta)] | y]$$

## Monotonicity / Chain rule

Additivity:  $g(\theta) = g_1(\theta) + g_2(\theta)$  if  $p(x|\theta) = p(x_1|\theta)p(x_2|\theta)$

Concavity:  $g_\lambda(\theta) \leq \lambda g_1(\theta) + (1-\lambda)g_2(\theta)$

## Cramer Rao inequality

unbiased estimator  $\mathbb{E}_\theta [\hat{\theta}(x)] = \theta \quad \forall \theta$

$$V_\theta = \mathbb{E}_\theta [(\hat{\theta}^i(x) - \theta^i)(\hat{\theta}^j(x) - \theta^j)]$$

$$V_\theta - \bar{g}^{-1}(\theta) \geq 0 \quad \text{positive semi-definite}$$