

# Local KL Geometry

The local KL geometry of statistical manifolds deviates from euclidean geometry in two ways: its quadratic part is curved and it contains data beyond its quadratic part. We examine how the new geometry distorts our euclidean picture of a large- $N$  sample as a tight Gaussian on an inner product space. We'll focus on the open simplex  $\mathcal{M}$  with  $D$  vertices.

The divergence  $KL : \mathcal{M} \times \mathcal{M} \rightarrow [0, \infty]$  given by

$$KL(\mathbf{q} : \mathbf{p}) = \mathbb{E}_{\mathbf{x} \sim \mathbf{q}}[\log(\mathbf{q}_x/\mathbf{p}_x)]$$

is smooth and vanishes on and only on the diagonal.

Let's first examine the chance  $\mathbb{P}[\hat{\mathbf{p}}; \mathbf{p}]$  of  $N$  i.i.d. samples from  $\mathbf{p}$  yielding the empirical distribution  $\hat{\mathbf{p}}$ . Fixing  $D$  and allowing proportionality constants to depend on  $N$ , we consider  $N$  large with respect to  $\hat{\mathbf{p}}$ ; that is, we take the sample granularity  $\alpha = \max_x 1/(\hat{\mathbf{p}}_x N)$  toward 0. Plugging Stirling's formula (here  $\mathbf{u}$  is uniform)

$$\binom{N}{\hat{\mathbf{p}}N} \propto \exp(NH(\hat{\mathbf{p}}) - (D/2)KL(\mathbf{u} : \hat{\mathbf{p}}) + O(\alpha))$$

into the lovely and routine formula (think of relative entropy)

$$\mathbb{P}[\hat{\mathbf{p}}; \mathbf{p}] \propto \binom{N}{\hat{\mathbf{p}}N} \exp(-NKL(\hat{\mathbf{p}} : \mathbf{p}) - NH(\hat{\mathbf{p}}))$$

we see

$$\mathbb{P}[\hat{\mathbf{p}}; \mathbf{p}] \propto \exp(-NKL(\hat{\mathbf{p}} : \mathbf{p}) - (D/2)KL(\mathbf{u} : \hat{\mathbf{p}}) + O(\alpha))$$

Now let's rescale per the central limit theorem. With  $KL(\mathbf{p} + \mathbf{v} : \mathbf{p}) = H(\mathbf{v}, \mathbf{v})/2 + J(\mathbf{v}, \mathbf{v}, \mathbf{v})/6 + o(v^3)$  and  $KL(\mathbf{u} : \mathbf{p} + \mathbf{v}) = KL(\mathbf{u} : \mathbf{p}) + \tilde{G}(\mathbf{v}) + o(v^1)$  we have for  $\mathbf{v} = \mathbf{u}/\sqrt{N}$  and  $\mathbf{u}$  bounded (so that  $\alpha \in O(1/N)$ ) that

$$\mathbb{P}[\mathbf{p} + \mathbf{u}/\sqrt{N}; \mathbf{p}] \propto \exp(-NH(\mathbf{u}, \mathbf{u})/2 - NJ(\mathbf{u}, \mathbf{u}, \mathbf{u})/6\sqrt{N} - D\tilde{G}(\mathbf{u})/2\sqrt{N} + O(1/N))$$

The  $\tilde{G}$  term translates the Gaussian toward  $\mathbf{u}$ ; the  $J$  term skews it (to have heavier tails) toward close vertices. For a sense of scale, note that  $\hat{\mathbf{p}}$ 's true distribution has mean  $\mathbf{p}$ , but this approximation suggests a shift on the order of  $\mathbf{u} \rightsquigarrow \mathbf{u} + 1/\sqrt{N}$  or  $\mathbf{v} \rightsquigarrow \mathbf{v} + 1/N$ . This is comparable to the resolution of  $\hat{\mathbf{p}}$ 's support and thus not a contradiction.

Let's visualize this for  $D = 3, N = 30, \mathbf{p} = (2/3, 3/12, 1/12)$ .

## Concentration in Graphs

**Markov.** Suppose we have a width- $W$  depth- $D$  'dense network' where each of  $D$  layers has  $W$  bit-valued nodes and where adjacent nodes have  $\leq \delta$  mutual information. How

does the total number of ‘on’ bits tend to deviate from the mean? (We might assume  $\mu = WD/2$  and call the deviation  $\Delta$ ).

First, if  $\delta = 0$ , then we have  $D$  independent subgaussians with parameter at most  $W$  so we expect concentration like  $-\Delta^2/2W^2D$ . A typical value for  $\Delta$  is at most  $W\sqrt{D}$ . At another extreme, if  $\delta = 1$ , then each layer can communicate its majority to the next, so we get no concentration.

Intuitively, if  $\delta$  mutual information corresponds to a probability  $p$  of corruption, then concentration should be as good as if we had  $pD$  many layers. Hmm... what’s a lower bound for  $p$  in terms of  $\delta \leq p \log(1/p) + p$ ? Well,  $\log(1/p) \leq \sqrt{1/p}$ , so  $\delta \leq 2\sqrt{p}$ .

$$p \log(1/p) \approx \delta \quad \delta / \log(1/\delta) \approx p \leq \delta \quad p \approx (\delta / \log(1/\delta)) / \log(\log(1/\delta)/\delta)$$