

Part I: Information Geometry

Color code

Definition

Theorem

Key Problem

1. The Center Piece of Information Theory

Definition: K-L Divergence (Kullback-Liebler)

For P, Q both probability distributions on the same finite alphabet \mathcal{X}

$$D(P||Q) \triangleq \sum_{x \in \mathcal{X}} P(x) \log \frac{P(x)}{Q(x)}$$

- Entropy and mutual information are both special cases

$$H(x) = H(P_x) = H(U) - D(P_x||U)$$
$$I(x; y) = D(P_{xy}||P_x P_y)$$

- **Information inequality**

$$D(P||Q) \geq 0, \quad \text{equality iff } P = Q$$

- **Convexity**

$$D(P||Q) \text{ is convex in } (P, Q)$$

- **Continuity**

Why K-L divergence matters?

Notation:

- Empirical distribution

$$\check{P}_x(a; x_1^n) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{(x_i=a)}$$

- Type class:

$$T_Q \triangleq \{x_1^n \in \mathcal{X}^n : \check{P}_x(\cdot; x_1^n) = Q(\cdot)\}$$

- Empirical Average

$$\frac{1}{n} \sum_{i=1}^n g(x_i) = \mathbb{E}_{x \sim \check{P}_x(\cdot; x_1^n)}[g(x)]$$

Sanov's Theorem

1. Probability of type class

$$\mathbb{P}_{x_1^n \sim \text{i.i.d. } P} (x_1^n \in T_Q) \doteq e^{-nD(Q||P)}$$

2. One type dominates: E is a subset of distributions on \mathcal{X} ,

$$\mathbb{P}_{x_1^n \sim \text{i.i.d. } P} (x_1^n \in \cup_{Q \in E} T_Q) \doteq e^{-n \cdot D(Q^*||P)}$$

where $Q^* = \arg \min_{Q \in E} D(Q||P)$

Quick review of the Channel Coding Theorem

- Transmit a codeword x_1^n , receive y_1^n ,
 (x_1^n, y_1^n) jointly typical w.r.t. P_{xy}
- Have some other "incorrect" codewords: $\tilde{x}_1^n[j]$, $j = 1, \dots, M$
each j : $(\tilde{x}_1^n[j], y_1^n) \sim P_x P_y$
- It's unlikely for an incorrect codeword to appear typical with the received
 $\mathbb{P}_{(\tilde{x}_1^n, y_1^n) \sim P_x P_y} ((\tilde{x}_1^n, y_1^n) \in T_{P_{xy}}) \doteq e^{-n \cdot D(P_{xy}||P_x P_y)} = e^{-n \cdot I(x;y)}$
- With $M = e^{nR}$ incorrect codewords, $R < I(x;y)$, the union bound of the above is still small.

Similar stories in rate distortion, error exponents,

2. Distance and Projection

- K-L divergence is a measure of distance between distributions
- There are many other ways to define divergence
eg. $f(\cdot)$ convex, continuous, and $f(1) = 0$

$$D_f(P||Q) = \sum_x Q(x) \cdot f\left(\frac{P(x)}{Q(x)}\right)$$

i-Projection, the binary hypothesis testing story

Consider x_1, \dots, x_n i.i.d. distributed from either P_0 or P_1 .

- Log-likelihood ratio test

$$\frac{1}{n} \sum_{i=1}^n \log \frac{P_1(x_i)}{P_0(x_i)} \underset{\hat{H}_0}{\overset{\hat{H}_1}{\geq}} \gamma$$

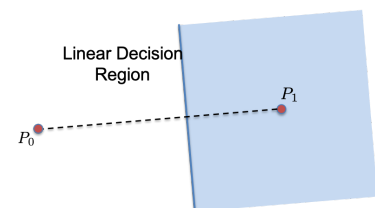
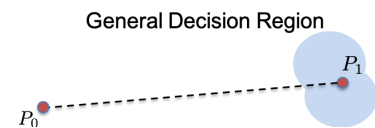
- The statistic is an empirical average

$$\frac{1}{n} \sum_{i=1}^n \log \frac{P_1(x_i)}{P_0(x_i)} = \mathbb{E}_{\mathbf{x} \sim \check{P}(\cdot; x_1^n)} \left[\log \frac{P_1}{P_0}(\mathbf{x}) \right]$$

- The decision region is a subset of *type classes*

$$E \triangleq \{Q : \mathbb{E}_{\mathbf{x} \sim Q} \left[\log \frac{P_1}{P_0}(\mathbf{x}) \right] \geq \gamma\}$$

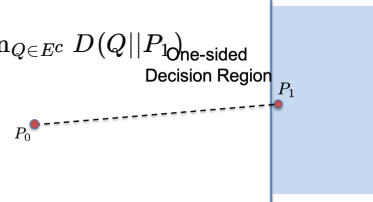
$$\text{claim } \hat{H}_1 \text{ iff } \mathbf{x}_1^n = x_1^n \in \bigcup_{Q \in E} T_Q$$



- Probability of Error

$$\mathbb{P}(H_0 \rightarrow \hat{H}_1) = \mathbb{P} \left(\mathbf{x}_1^n \in \bigcup_{Q \in E} T_Q \mid H_0 \right) \doteq e^{-n \cdot \min_{Q \in E} D(Q \| P_0)}$$

$$\mathbb{P}(H_1 \rightarrow \hat{H}_0) = \mathbb{P} \left(\mathbf{x}_1^n \in \bigcup_{Q \in E^c} T_Q \mid H_1 \right) \doteq e^{-n \cdot \min_{Q \in E^c} D(Q \| P_1)}$$



- The optimization problem

$$Q^* = \arg \min_{Q: \mathbb{E}_Q[f(\mathbf{x})] > \gamma} D(Q \| P_0)$$

- Dominating error event $Q_{0 \rightarrow 1, \gamma}^*$
- testing statistic $f(\mathbf{x}) = \log P_1(\mathbf{x}) / P_0(\mathbf{x})$

Definition: Exponential Family (1-D)

$$\mathcal{E}(P_0, f) \triangleq \{P_t, t \in [0, 1] : P_t(x) = P_0(x) \cdot e^{t \cdot f(x) - \alpha(t)}, \forall x\}$$

- P_0 : a starting point
- $f(\cdot)$: natural statistic (meaning later)
- $\alpha(t)$: normalization factor

$$e^{\alpha(t)} = \sum_x P_0(x) \cdot e^{t \cdot f(x)} = \mathbb{E}_{\mathbf{x} \sim P_0} [e^{t \cdot f(\mathbf{x})}]$$

also called the log-moment generation function.

- viewed as "exponential tilting" on P_0 according to $f(\cdot)$.

- Empirical average

$$\eta(t) \triangleq \mathbb{E}_{\mathbf{x} \sim P_t} [f(\mathbf{x})]$$

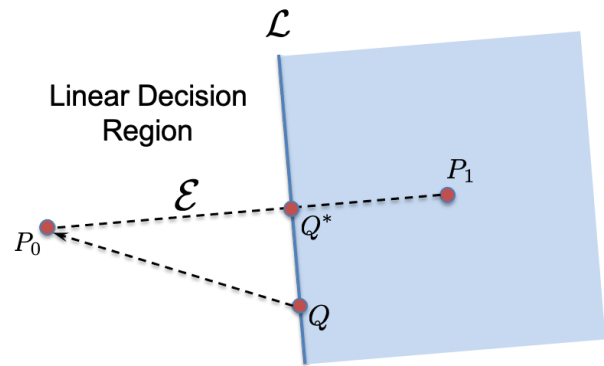
- A number of nice properties

$$\frac{\partial^2}{\partial t^2} D(P_t || P_0) = \frac{\partial}{\partial t} \eta(t) = \text{var}_{\mathbf{x} \sim P_t} [f(\mathbf{x})] = \mathcal{I}_t$$

- $\eta(t)$ monotonically increase with t
- Connection between Fisher information \mathcal{I}_t and K-L divergence

Definition: Linear Family

$$\mathcal{L}(f, \gamma) \triangleq \{Q : \mathbb{E}_{\mathbf{x} \sim Q} [f(\mathbf{x})] = \gamma\}$$



Theorem: Pythagorean

$$\forall Q \in \mathcal{L}(f, \gamma) :$$

$$D(Q || P_0) = D(Q || Q^*) + D(Q^* || P_0)$$

$$\text{where } Q^* \in \mathcal{L}(f, \gamma) \cap \mathcal{E}(P_0, f)$$

- Unique intersection since $\eta(t) \triangleq \mathbb{E}_{\mathbf{x} \sim P_t} [f(\mathbf{x})]$ monotonic increase with t .
 $Q^* = P_{t^*} \in \mathcal{E}$, with $\mathbb{E}_{\mathbf{x} \sim P_{t^*}} [f(\mathbf{x})] = \gamma$:

$$\begin{aligned}
D(Q||Q^*) &= \mathbb{E}_{x \sim Q} \left[\log \frac{Q(x)}{Q^*(x)} \right] = \mathbb{E}_{x \sim Q} \left[\log \frac{Q(x)}{P_{t^*}(x)} \right] \\
&= \mathbb{E}_{x \sim Q} \left[\log \frac{Q(x)}{P_0(x) \cdot e^{t^* \cdot f(x) - \alpha(t^*)}} \right] \\
&= \mathbb{E}_{x \sim Q} \left[\log \frac{Q(x)}{P_0(x)} \right] - \mathbb{E}_{x \sim Q} [t^* f(x) - \alpha^*(t)] \quad \text{C} \\
&= \mathbb{E}_{x \sim Q} \left[\log \frac{Q(x)}{P_0(x)} \right] - \mathbb{E}_{x \sim Q^*} [t^* f(x) - \alpha^*(t)] \\
&= \mathbb{E}_{x \sim Q} \left[\log \frac{Q(x)}{P_0(x)} \right] - \mathbb{E}_{x \sim Q^*} \left[\log \frac{P_0(x) \cdot e^{t^* f(x) - \alpha^*(t)}}{P_0(x)} \right] \\
&= D(Q||P_0) - D(Q^*||P_0)
\end{aligned}$$

Corollary: Typical Error Event Occurs on Exponential Family

$$Q^* = \arg \min_{Q: \mathbb{E}_Q[f(x)] > \gamma} D(Q||P_0)$$

has $Q^* \in \mathcal{E}(P_0, f)$, with $f(x) = \log \left(\frac{P_1(x)}{p_0(x)} \right)$ and $\mathbb{E}_{x \sim Q^*} [f(x)] = \gamma$.

Definition: Q^* is called the i-projection of P_0 to the linear family $\mathcal{L}(f, \gamma)$.

Takeaway message:

- Hypothesis testing is about operations on the empirical distribution, in functional space;
- Each problem has a pair P_0, P_1 , and the exponential family associated;
- Projection of the observed empirical distribution to the exponential family.

m-Projection: the Learning Story

Suppose we observe some data samples x_1^n with empirical distribution $\check{P}(\cdot; x_1^n)$.

We know that the true model belongs to a parameterized family

$$\mathcal{P} \triangleq \{P(\cdot; \theta); \theta \in \mathbb{R}\}$$

often chosen as an exponential family.

Maximum Likelihood estimate of the unknown parameter θ .

$$\hat{\theta}_{\text{ML}}(x^n) = \arg \max_{\hat{\theta}} \frac{1}{n} \sum_{i=1}^n \log P(x_i; \hat{\theta})$$

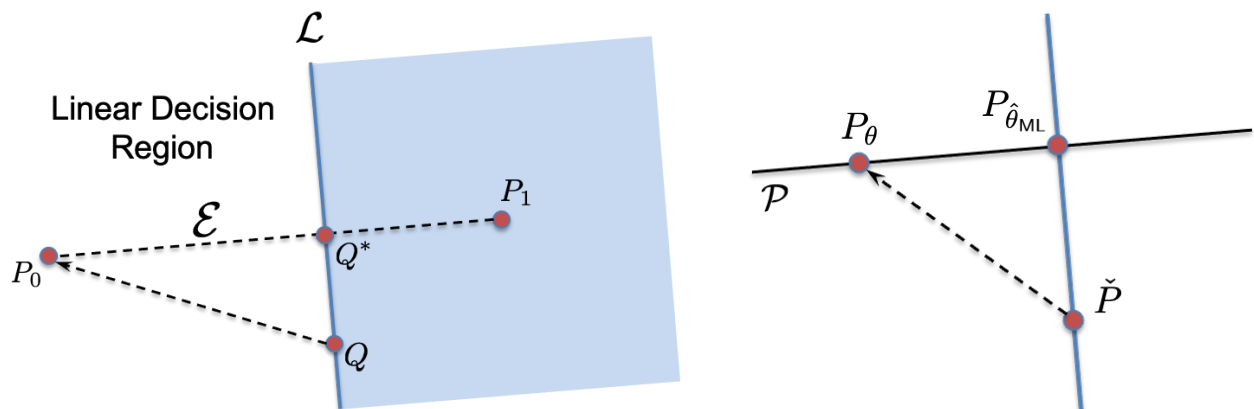
- Usually assume the family to be smooth, $\frac{\partial}{\partial \theta} P(x, \theta)$ exist, finite
- Can have higher dimensional parameters
- Why do maximum likelihood estimate?
- Distribution matching

$$\begin{aligned} \arg \max_{\hat{\theta}} \frac{1}{n} \sum_{i=1}^n \log P(x_i; \hat{\theta}) &= \arg \max_{\hat{\theta}} \mathbb{E}_{\mathbf{x} \sim \check{P}_{\mathbf{x}}(\cdot; x_1^n)} \left[\log P(\mathbf{x}; \hat{\theta}) \right] \\ &= \arg \min_{\hat{\theta}} \mathbb{E}_{\mathbf{x} \sim \check{P}(\cdot; x_1^n)} \left[\log \frac{\check{P}(\mathbf{x}; x_1^n)}{P(\mathbf{x}; \hat{\theta})} \right] \\ &= \arg \min_{\hat{\theta}} D \left(\check{P}(\cdot; x_1^n) \parallel P(\cdot; \hat{\theta}) \right) \end{aligned}$$

Definition: $P(\cdot; \hat{\theta}_{\text{ML}})$ is called the m-projection of \check{P} to the model family \mathcal{P} .

Corollary: if \mathcal{P} is an exponential family, $\mathcal{P} = \mathcal{E}(P_0, f)$, and suppose $\mathbb{E}_{\mathbf{x} \sim \check{P}}[f(\mathbf{x})] = \gamma$, then

$$P(\cdot; \hat{\theta}_{\text{ML}}) \in \mathcal{P} \cap \mathcal{L}(f, \gamma)$$



Takeaway message:

- A model family is a manifold/plane in the space of distributions;
- Learning is also a projection, from the observed empirical distribution to the model family.

3. The Geometry of Information Theory and Learning

- A number of information theory results presented as geometric stories:
 - Rate Distortion, "Rate-distortion theory: A mathematical basis for data compression, Englewood Cliffs, NJ: Prentice-Hall, 1971."
 - Error Exponent
 - Csiszar's book

now publishers - Information Theory and Statistics: A Tutorial
Publication Date: 15 Dec 2004 Download extract Abstract This tutorial is concerned with applications of information theory concepts in statistics, in the finite alphabet setting. The information

<https://www.nowpublishers.com/article/Details/CIT-004>

Information Theory
and Statistics:
A Tutorial
Imre Csiszár
Peter C. Shields

- What is difficult about this?
 - The geometry is complex.

Shun'ichi Amari - Wikipedia

Shun'ichi Amari, is a Japanese scholar born in 1936 in Tokyo, Japan. He majored in Mathematical Engineering in 1958 from the University of Tokyo then graduated in

[W https://en.wikipedia.org/wiki/Shun%27ichi_Amari](https://en.wikipedia.org/wiki/Shun%27ichi_Amari)



- Fisher information

$$\frac{\partial^2}{\partial t^2} D(P_t || P_0) = \frac{\partial}{\partial t} \eta(t) = \text{var}_{\mathbf{x} \sim P_t} [f(\mathbf{x})] = \mathcal{I}_t$$

but $\mathcal{I}_t \geq 0$ can be an arbitrary function of t .

So $D(P_t||P_0)$ is a convex function of t , but not clear how convex.

- If you have learned Cramer-Rao bound ...
- What we need is a lot more.
 - Broadcast channels: $P_{y|x}, P_{z|x}$. Even if $I(x;y) > I(x;z)$, doesn't mean the channel $x \rightarrow z$ is degraded.

Dependence is not a single dimensional concept.

- Mismatched detection, universal detection: what happens if we didn't use the right $f(x) = \log \frac{P_1(x)}{P_0(x)}$ to make decision, but used a different $f'(\cdot)$?

How bad are imperfect statistic models?

- Increasing the dimensionality of \mathcal{E} , what collection/sequence of statistics

$$P_x(x; \underline{\theta}) = P_0(x) \cdot \exp \left[\sum_{i=1}^k \theta_i \cdot f_i(x) - \alpha(\underline{\theta}) \right]$$

What statistic is more valuable in learning?

- What happens with each iteration and each mini-batch of samples?

Evolution and convergence of learned models in functional space.

- From input/output neural networks to Transfer Learning, Multi-Modal Learning

Network information theory and more complex learning tasks.

- There are often too many distributions to worry about
 - The ground truth
 - The parameterized family
 - The empiricals
 - The current model and the updates
 - Restrictions, side information, loss
 - Tuning of design parameters

- **Basically:** we cannot write it very clean for 1-D problems with 2 distributions, but we need some analysis for multi-dimensional problems with many distributions.

What is Geometry and Why Geometry?

- Distance \longrightarrow inner product, projection, basis, coordinates (Hilbert Space for distributions)
- Space of functions and Space of distributions.