

1. Properties of f -Divergences.

For any $P, Q \in \mathcal{P}(\mathcal{X})$ probability measures on the same probability space, dominated by a common measure $P, Q \ll \lambda$, recall that

$$D_f(P\|Q) := \mathbb{E}_Q \left[f \left(\frac{dP/d\lambda}{dQ/d\lambda} \right) \right]$$

where f is a convex function satisfying the assumption given in class and $dP/d\lambda$ is the Radon-Nikodym derivative of P with respect to λ . Prove the following properties:

- (a) Non-Negativity: $D_f(P\|Q) \geq 0$ with equality if and only if $P = Q$.

Solution. By the definition of f -Divergence we have

$$D_f(P\|Q) = \mathbb{E}_Q \left[f \left(\frac{dP/d\lambda}{dQ/d\lambda} \right) \right] \geq f \left(\mathbb{E}_Q \left[\frac{dP/d\lambda}{dQ/d\lambda} \right] \right) \geq f(1) = 0$$

where the first inequality follows from Jensen's inequality and the second from convexity of f . To prove equality, first assume that $P = Q$. Then,

$$D_f(P\|Q) = \mathbb{E}_Q[f(1)] = 0.$$

Now assume $D_f(P\|Q) = 0$. Then $f \left(\frac{dP/d\lambda}{dQ/d\lambda} \right) = 0 \implies \frac{dP/d\lambda}{dQ/d\lambda} = 1$ since f is strongly convex at 1. It follows that $P = Q$. ■

- (b) Joint Convexity: The map $(P, Q) \mapsto D_f(P\|Q)$ is (jointly) convex.

Solution. For any function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ we can define the perspective of f as the function $g : \mathbb{R}^n \times \mathbb{R}_{>0} \rightarrow \mathbb{R}$ such that

$$g(x, y) = yf\left(\frac{x}{y}\right), \quad \text{dom } g = \left\{ (x, y) \mid \frac{x}{y} \in \text{dom } f, y > 0 \right\}$$

For a convex function f , the perspective of f is also convex. That is,

$$\begin{aligned} g(\alpha(x_1, y_1) + (1 - \alpha)(x_2, y_2)) &\leq \alpha g(x_1, y_1) + (1 - \alpha)g(x_2, y_2) \\ &= \alpha y_1 f\left(\frac{x_1}{y_1}\right) + (1 - \alpha)y_2 f\left(\frac{x_2}{y_2}\right) \end{aligned}$$

for all $\alpha \in [0, 1]$ and each (x_i, y_i) in the domain of g . Define $\tilde{D}_f(P, Q) = D_f(P\|Q)$.

Then for $(P_1, Q_1), (P_2, Q_2) \in \mathcal{P}(\mathcal{X} \times \mathcal{X})$ such that $P_1, P_2, Q_1, Q_2 \ll \lambda$

$$\begin{aligned} \tilde{D}_f(\alpha(P_1, Q_1) + (1 - \alpha)(P_2, Q_2)) &= \int_{\mathcal{X}} f\left(\frac{\alpha \frac{dP_1}{d\lambda}(x) + (1 - \alpha) \frac{dP_2}{d\lambda}(x)}{\alpha \frac{dQ_1}{d\lambda}(x) + (1 - \alpha) \frac{dQ_2}{d\lambda}(x)}\right) d\lambda \\ &\leq \int_{\mathcal{X}} \alpha f\left(\frac{\frac{dP_1}{d\lambda}(x)}{\frac{dQ_1}{d\lambda}(x)}\right) d\lambda + (1 - \alpha) \int_{\mathcal{X}} f\left(\frac{\frac{dP_2}{d\lambda}(x)}{\frac{dQ_2}{d\lambda}(x)}\right) d\lambda \\ &= \alpha \tilde{D}_f(P_1, Q_1) + (1 - \alpha) \tilde{D}_f(P_2, Q_2) \end{aligned}$$

Thus the mapping $(P, Q) \mapsto D_f(P\|Q)$ is convex. ■

- (c) Conditioning Increases f -Divergences: For $P_X \in \mathcal{P}(\mathcal{X})$ and two transition kernels $P_{Y|X}$ and $Q_{Y|X}$ from \mathcal{X} to \mathcal{Y} , consider the probability measures $P_{XY} := P_X P_{Y|X}$ and $Q_{XY} := Q_X P_{Y|X}$ on $\mathcal{X} \times \mathcal{Y}$. Denoting P_Y and Q_Y as their marginals of \mathcal{Y} , show that

$$D_f(P_Y\|Q_Y) \leq D_f(P_{Y|X}\|Q_{Y|X}|P_X)$$

Solution.

$$\begin{aligned} D_f(P_Y\|Q_Y) &= D_f\left(\mathbb{E}_{P_X}[P_{Y|X}(\cdot|X)]\|\mathbb{E}_{P_X}[Q_{Y|X}(\cdot|X)]\right) \\ &\leq \mathbb{E}_{P_X}\left[D_f(P_{Y|X}(\cdot|x)\|Q_{Y|X}(\cdot|x))\right] \quad (\text{Jensen's Inequality}) \\ &= D_f(P_{Y|X}\|Q_{Y|X}|P_X) \end{aligned}$$

■

- (d) Joint vs. Marginal: For $P_X, Q_X \in \mathcal{P}(\mathcal{X})$ and a transition kernel $P_{Y|X}$, define $P_{XY} := P_X P_{Y|X}$ and $Q_{XY} := Q_X P_{Y|X}$ on $\mathcal{X} \times \mathcal{Y}$. Show that

$$D_f(P_X\|Q_X) = D_f(P_{XY}\|Q_{XY})$$

Solution.

$$\begin{aligned}
 D_f(P_{XY} \| Q_{XY}) &= \mathbb{E}_{Q_{XY}} \left[f \left(\frac{dP_{XY}}{dQ_{XY}} \right) \right] = \mathbb{E}_{Q_{XY}} \left[f \left(\frac{dP_X P_{Y|X}}{dQ_X P_{Y|X}} \right) \right] \\
 &= \int_{\mathcal{X}} \int_{\mathcal{Y}} f \left(\frac{dP_X P_{Y|X}}{dQ_X P_{Y|X}} \right) dQ_{X,Y}(x, y) = \int_{\mathcal{X}} f \left(\frac{dP_X}{dQ_X} \right) dQ_X(x) \\
 &= \mathbb{E}_{Q_X} \left[f \left(\frac{dP_X}{dQ_X} \right) \right] = D_f(P_X \| Q_X)
 \end{aligned}$$

■

2. Example of Data Processing Inequality.

Let $(\mathcal{X}, \mathcal{F})$ be a measurable space. Use the Data Processing Inequality to show that for any two measurable P, Q on $(\mathcal{X}, \mathcal{F})$ and any $A \in \mathcal{F}$:

$$D_f(P \| Q) \geq \sup_{A \in \mathcal{F}} \left\{ (1 - Q(A)) f \left(\frac{1 - P(A)}{1 - Q(A)} \right) + Q(A) f \left(\frac{P(A)}{Q(A)} \right) \right\}$$

Solution. Let $\mathbb{1}^A(\cdot | \cdot)$ denote the Dirac measure, where

$$\mathbb{1}^A(1|x) = \begin{cases} 1, & x \in A \\ 0, & \text{otherwise} \end{cases}$$

By taking the expectation with respect to P_X and Q_X , we get measures P_Y, Q_Y that correspond to the Bernoulli variable $Ber(P(A))$ and $Ber(Q(A))$ respectively. Note that this would be a discrete random variable with sample space $\{0, 1\}$ that's absolutely continuous with respect to the counting measure. Computation of the divergence yields

$$D_f(P_Y \| Q_Y) = (1 - Q(A)) f \left(\frac{1 - P(A)}{1 - Q(A)} \right) + Q(A) f \left(\frac{P(A)}{Q(A)} \right).$$

Using the data processing inequality, we get

$$D_f((\|P), Q) \geq D_f(P_Y \| Q_Y) \tag{1}$$

$$= (1 - Q(A)) f \left(\frac{1 - P(A)}{1 - Q(A)} \right) + Q(A) f \left(\frac{P(A)}{Q(A)} \right) \quad \forall A \in \mathcal{F}. \tag{2}$$

Since the inequality holds $\forall A \in \mathcal{X}$ we arrive at the conclusion that $D_f((\|P), Q)$ is an upper bound to

$$\left\{ (1 - Q(A)) f \left(\frac{1 - P(A)}{1 - Q(A)} \right) + Q(A) f \left(\frac{P(A)}{Q(A)} \right) \mid A \in \mathcal{F} \right\},$$

and by definition of sup, we get

$$D_f(P\|Q) \geq \sup_{A \in \mathcal{F}} \left\{ (1 - Q(A))f\left(\frac{1 - P(A)}{1 - Q(A)}\right) + Q(A)f\left(\frac{P(A)}{Q(A)}\right) \right\}$$

■

3. f -Divergences, Metrics, and Mismatched Supports.

For the KL divergence $D_{\text{KL}}(\cdot\|\cdot)$ and $\chi^2(\cdot\|\cdot)$ as shown in class, show that:

(a) $\delta_{\text{TV}}(\cdot, \cdot)$ is a metric on $\mathcal{P}(\mathcal{X})$.

Solution.

i. Identity:

$$\delta_{\text{TV}}(P, P) = \mathbb{E}_Q \left[\frac{1}{2} \left| \frac{dP/d\lambda}{dP/d\lambda} - 1 \right| \right] = 0$$

ii. Symmetry:

$$\begin{aligned} \delta_{\text{TV}}(P, Q) &= \mathbb{E}_Q \left[\frac{1}{2} \left| \frac{dP/d\lambda}{dQ/d\lambda} - 1 \right| \right] = \int_{\mathcal{X}} \left| \frac{dP}{d\lambda}(x) - \frac{dQ}{d\lambda}(x) \right| d\lambda \\ &= \int_{\mathcal{X}} \left| \frac{dQ}{d\lambda}(x) - \frac{dP}{d\lambda}(x) \right| d\lambda \\ &= \mathbb{E}_P \left[\frac{1}{2} \left| \frac{dQ/d\lambda}{dP/d\lambda} - 1 \right| \right] \\ &= \delta_{\text{TV}}(Q, P) \end{aligned}$$

iii. Triangle Inequality: For $P, Q, R \in \mathcal{P}(\mathcal{X})$,

$$\begin{aligned} \delta_{\text{TV}}(P, Q) + \delta_{\text{TV}}(Q, R) &= \int_{\mathcal{X}} \left| \frac{dP}{d\lambda}(x) - \frac{dQ}{d\lambda}(x) \right| d\lambda + \int_{\mathcal{X}} \left| \frac{dQ}{d\lambda}(x) - \frac{dR}{d\lambda}(x) \right| d\lambda \\ &= \int_{\mathcal{X}} \left| \frac{dP}{d\lambda}(x) - \frac{dQ}{d\lambda}(x) \right| + \left| \frac{dQ}{d\lambda}(x) - \frac{dR}{d\lambda}(x) \right| d\lambda \\ &\geq \int_{\mathcal{X}} \left| \frac{dP}{d\lambda}(x) - \frac{dR}{d\lambda}(x) \right| d\lambda \\ &= \delta_{\text{TV}}(P, R) \end{aligned}$$

■

(b) $D_{\text{KL}}(P\|Q) = \chi^2(P\|Q) = \infty$ whenever $P \not\ll Q$.

Solution.

$$\chi^2(P\|Q) = \int_{\mathcal{X}} \left(\left(\frac{dP/d\lambda}{dQ/d\lambda} \right)^2 - 1 \right) d\lambda$$

which will blow up at values where Q attains 0 and P does not which will happen since $P \not\ll Q$.

$$D_{\text{KL}}(P\|Q) = \int_{\mathcal{X}} \frac{dP}{d\lambda}(x) \log \left(\frac{dP/d\lambda}{dQ/d\lambda} \right) d\lambda$$

will blow up if $P \not\ll Q$ in a similar fashion since the denominator in the log will be zero when P is not zero which prevents our convention of $0f(0/0) = 0$ yielding an infinite divergence. ■

- (c) $\delta_{\text{TV}}(P, Q)$ attains its maximal value of 1 when $\text{supp}(P) \cap \text{supp}(Q) = \emptyset$.

Solution.

$$\begin{aligned} \int_{\mathcal{X}} \frac{1}{2} \left| \frac{dP}{d\lambda}(x) - \frac{dQ}{d\lambda}(x) \right| d\lambda &\leq \int_{\mathcal{X}} \frac{1}{2} \left| \frac{dP}{d\lambda}(x) \right| + \left| \frac{dQ}{d\lambda}(x) \right| d\lambda \leq \int_{\mathcal{X}} \frac{1}{2} \left| \frac{dP}{d\lambda}(x) \right| d\lambda + \int_{\mathcal{X}} \left| \frac{dQ}{d\lambda}(x) \right| d\lambda \\ &\leq \frac{1}{2} + \frac{1}{2} = 1 \end{aligned}$$

with equality if and only if $\text{supp}(P) \cap \text{supp}(Q) = \emptyset$. ■

- (d) Explain why the previous property is undesired when performing generative modeling $\inf_{\theta \in \Theta} \delta_{\text{TV}}(P, Q_{\theta})$ of a data distribution P via a parametrized family $\{Q_{\theta}\}_{\theta \in \Theta}$ under divergence δ .

Solution.

When the supports are disjoint we have a constant distance (which is at its max) between any distribution Q_{θ} and our data distribution P . This means the model will have the same error throughout training preventing the model from learning what works and what does not in its generations. ■

4. Jensen-Shannon Divergence

Let $f(x) = x \log \left(\frac{2x}{x+1} \right) + \log \left(\frac{2}{x+1} \right)$. Show that:

- (a) $f : (0, \infty) \rightarrow \mathbb{R}$ is a convex function, with $f(1) = 0$, which is strictly convex around 1.

Solution. We compute the second derivative to be

$$\begin{aligned} f'(x) &= \log \left(\frac{2x}{x+1} \right) \\ f''(x) &= \frac{1}{x^2 + x} \end{aligned}$$

which reveals $f''(x) > 0$ for $x > 0$ giving that f is convex. Evaluating $f(1)$ gives

$$f(1) = 1 \cdot \log\left(\frac{2}{2}\right) + \log\left(\frac{2}{2}\right) = 0 + 0 = 0.$$

as desired. Finally, to see that f is strictly convex around 1 observe that the second derivative is positive at all points in our domain. This implies that f is strictly convex. The minimum of f is found to be $x = 1$ after analyzing the first derivative. Thus f is strictly convex around 1. ■

(b) Let $\text{JSD}(P\|Q)$ be the f -divergence induced by the above f . Prove that

$$\text{i. } \text{JSD}(P\|Q) = D_{\text{KL}}\left(P\|\frac{P+Q}{2}\right) + D_{\text{KL}}\left(Q\|\frac{P+Q}{2}\right)$$

Solution. We iron both expressions out and realize they are the same.

$$\begin{aligned} \text{JSD}(P\|Q) &= \mathbb{E}_Q \left[\frac{dP/d\lambda}{dQ/d\lambda} \log \left(\frac{2 \frac{dP/d\lambda}{dQ/d\lambda}}{\frac{dP/d\lambda}{dQ/d\lambda} + 1} \right) + \log \left(\frac{2}{\frac{dP/d\lambda}{dQ/d\lambda} + 1} \right) \right] \\ &= \int_{\mathcal{X}} \left[\frac{dP/d\lambda}{dQ/d\lambda} \log \left(\frac{2 \frac{dP/d\lambda}{dQ/d\lambda}}{\frac{dP/d\lambda}{dQ/d\lambda} + 1} \right) + \log \left(\frac{2}{\frac{dP/d\lambda}{dQ/d\lambda} + 1} \right) \right] dQ \\ &= \int_{\mathcal{X}} \frac{dP}{d\lambda} \log \left(\frac{2dP/d\lambda}{dP/d\lambda + dQ/d\lambda} \right) d\lambda + \int_{\mathcal{X}} \frac{dQ}{d\lambda} \log \left(\frac{2dQ/d\lambda}{dP/d\lambda + dQ/d\lambda} \right) d\lambda \end{aligned}$$

$$\begin{aligned}
D_{\text{KL}}\left(P\left\|\frac{P+Q}{2}\right.\right) + D_{\text{KL}}\left(Q\left\|\frac{P+Q}{2}\right.\right) &= \mathbb{E}_{\frac{P+Q}{2}} \left[\frac{dP/d\lambda}{d\left(\frac{P+Q}{2}\right)/d\lambda} \log \left(\frac{dP/d\lambda}{d\left(\frac{P+Q}{2}\right)/d\lambda} \right) \right] \\
&\quad + \mathbb{E}_{\frac{P+Q}{2}} \left[\frac{dQ/d\lambda}{d\left(\frac{P+Q}{2}\right)/d\lambda} \log \left(\frac{dQ/d\lambda}{d\left(\frac{P+Q}{2}\right)/d\lambda} \right) \right] \\
&= \int_{\mathcal{X}} \frac{dP/d\lambda}{d\left(\frac{P+Q}{2}\right)/d\lambda} \log \left(\frac{dP/d\lambda}{d\left(\frac{P+Q}{2}\right)/d\lambda} \right) d\left(\frac{P+Q}{2}\right) \\
&\quad + \int_{\mathcal{X}} \frac{dQ/d\lambda}{d\left(\frac{P+Q}{2}\right)/d\lambda} \log \left(\frac{dQ/d\lambda}{d\left(\frac{P+Q}{2}\right)/d\lambda} \right) d\left(\frac{P+Q}{2}\right) \\
&= \int_{\mathcal{X}} \frac{dP}{d\lambda} \log \left(\frac{dP/d\lambda}{d\left(\frac{P+Q}{2}\right)/d\lambda} \right) d\lambda + \int_{\mathcal{X}} \frac{dQ}{d\lambda} \log \left(\frac{dQ/d\lambda}{d\left(\frac{P+Q}{2}\right)/d\lambda} \right) d\lambda \\
&= \int_{\mathcal{X}} \frac{dP}{d\lambda} \log \left(\frac{2dP/d\lambda}{dP/d\lambda + dQ/d\lambda} \right) d\lambda + \int_{\mathcal{X}} \frac{dQ}{d\lambda} \log \left(\frac{dQ/d\lambda}{dP/d\lambda + dQ/d\lambda} \right) d\lambda
\end{aligned}$$

Thus the expressions are equivalent. ■

ii. JSD $(P\|Q)$ is maximized at $2\log 2$.

Solution. We note that

$$\begin{aligned}
D_{\text{KL}}\left(P\left\|\frac{P+Q}{2}\right.\right) &= \int_{\mathcal{X}} \frac{dP}{d\lambda} \log \left(\frac{2dP/d\lambda}{dP/d\lambda + dQ/d\lambda} \right) d\lambda \leq \int_{\mathcal{X}} \frac{dP}{d\lambda} \log \left(\frac{2dP/d\lambda}{dP/d\lambda + 0} \right) d\lambda \\
&= \log(2) \int_{\mathcal{X}} \frac{dP}{d\lambda} d\lambda = \log(2).
\end{aligned}$$

Note that by construction of the proof that equality only holds when

$$\text{supp}(P) \cap \text{supp}(Q) = \emptyset.$$

We similarly conclude that

$$D_{\text{KL}}\left(Q\left\|\frac{P+Q}{2}\right.\right) \leq \log(2).$$

Therefore

$$\text{JSD}(P\|Q) = D_{\text{KL}}\left(P\|\frac{P+Q}{2}\right) + D_{\text{KL}}\left(Q\|\frac{P+Q}{2}\right) \leq \log(2) + \log(2) = 2\log(2)$$

■

5. f -Divergences Variational Formula.

The convex conjugate of a function $f : I \rightarrow \mathbb{R}$ is $f^*(y) = \sup_{x \in I} yx - f(x)$. We saw the following variational representation of f -divergences:

$$D_f(P\|Q) = \sup_{g: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{E}_P[g(X)] - \mathbb{E}_Q[f^*(g(X))],$$

where the supremum is over all measurable g for which the expectations are finite. Show that:

(a) $D_f(P\|Q) \geq \sup_{g: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{E}_P[g(X)] - \mathbb{E}_Q[f^*(g(X))]$ when supremising over all g as above.

Solution. For $P, Q \in \mathcal{P}(\mathcal{X})$, $I^* = \{y \in \mathbb{R} \mid yx - f(x) < \infty\}$,

$$D_f(P\|Q) = \int_{x \in \mathcal{X}} f\left(\frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)}\right) dQ(x) = \int_{x \in \mathcal{X}} \sup_{y \in I^*} \left(\frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)} y - f^*(y)\right) dQ(x).$$

Pick any $g : \mathcal{X} \rightarrow I^*$ and set $y = g(x)$ to get a lower bound:

$$\begin{aligned} \int_{x \in \mathcal{X}} \sup_{y \in I^*} \left(\frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)} y - f^*(y)\right) dQ(x) &\geq \int_{x \in \mathcal{X}} \left(\frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)} g(x) - f^*(g(x))\right) dQ(x) \\ &= \int_{x \in \mathcal{X}} \frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)} g(x) dQ(x) - \int_{x \in \mathcal{X}} f^*(g(x)) dQ(x) \\ &= \int_{x \in \mathcal{X}} \frac{dP}{d\lambda}(x) g(x) d\lambda - \int_{x \in \mathcal{X}} f^*(g(x)) dQ(x) \\ &= \mathbb{E}_P[g(x)] - \mathbb{E}_Q[f^*(g(x))] \end{aligned}$$

Supremizing over all measurable $g : \mathcal{X} \rightarrow I^*$, we get

$$D_f(P\|Q) \geq \sup_{g: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{E}_P[g(x)] - \mathbb{E}_Q[f^*(g(x))]$$

as desired. ■

(b) Derive the following variational formulas by computing convex conjugates:

i. $D_{\text{KL}}(P\|Q) = 1 + \sup_{g:\mathcal{X} \rightarrow \mathbb{R}} \mathbb{E}_P[g(X)] - \mathbb{E}_Q[e^{g(X)}].$

Solution. For KL divergence we have

$$\begin{aligned} f(x) &= x \log(x), \\ f^*(y) &= \sup_{x \in (0, \infty)} xy - x \log(x). \end{aligned}$$

Letting $s(x) = xy - x \log(x)$, we compute the location of the maxima of $s(x)$ by

$$\begin{aligned} \left. \frac{ds(x)}{dx} \right|_{x>0} &= y - 1 - \log(x) \\ &\rightarrow x = e^{y-1}. \end{aligned}$$

Thus

$$\begin{aligned} f^*(y) &= ye^{y-1} - e^{y-1} \log(e^{y-1}) \\ &= ye^{y-1} - (y-1)e^{y-1} \\ &= e^{y-1}. \end{aligned}$$

It follows that

$$D_{\text{KL}}(P\|Q) = \sup_{g:\mathcal{X} \rightarrow \mathbb{R}} (\mathbb{E}_P[g(x)] - \mathbb{E}_Q[f^*(g(x))]) = \sup_{g:\mathcal{X} \rightarrow \mathbb{R}} (\mathbb{E}_P[g(x)] - \mathbb{E}_Q[e^{g(x)-1}]).$$

Letting $h(x) = g(x) - 1$ we get the desired expression

$$\begin{aligned} D_{\text{KL}}(P\|Q) &= \sup_{h:\mathcal{X} \rightarrow \mathbb{R}} (\mathbb{E}_P[h(x) + 1] - \mathbb{E}_Q[e^{h(x)}]) \\ &= 1 + \sup_{h:\mathcal{X} \rightarrow \mathbb{R}} (\mathbb{E}_P[h(x)] - \mathbb{E}_Q[e^{h(x)}]). \end{aligned}$$

■

ii. $\delta_{\text{TV}}(P, Q) = \sup_{\|g\|_{\infty} \leq 1} \frac{1}{2} (\mathbb{E}_P[g(X)] - \mathbb{E}_Q[g(X)]).$

Solution. For TV distance we have that

$$\begin{aligned} f(x) &= \frac{1}{2}|x - 1|, \\ f^*(y) &= \sup_{x \in (0, \infty)} xy - \frac{1}{2}|x - 1|. \end{aligned}$$

Letting $s(x) = xy - \frac{1}{2}|x - 1|$, we compute the location of the maxima of $s(x)$ by

$$\left. \frac{ds(x)}{dx} \right|_{x>0} = \begin{cases} y + \frac{1}{2}, & x \leq 1 \\ y - \frac{1}{2}, & x > 1 \end{cases}$$

Note that for $f(y)$ to be bounded, we need the derivative for when $x > 1$ to be negative as otherwise the function will explode to infinity. This means that we need $|y| \leq \frac{1}{2}$. And for when $x \leq 1$, the max is attained at $x = 1$ and it is bounded by y . So supremizing over x we have

$$f^*(y) = \begin{cases} y, & |y| \leq \frac{1}{2} \\ \infty, & \text{otherwise} \end{cases}$$

Thus, we have

$$\begin{aligned} \delta_{\text{TV}}(P, Q) &= \sup_{\|g\|_{\infty} \leq \frac{1}{2}} \mathbb{E}_P[g(x)] + \mathbb{E}_Q[f^*(g(x))] \\ &= \sup_{\|g\|_{\infty} \leq 1} \frac{1}{2} (\mathbb{E}_P[g(x)] + \mathbb{E}_Q[g(x)]) \end{aligned}$$

■

iii. $\chi^2(P||Q) = \sup_{g:\mathcal{X} \rightarrow \mathbb{R}} \mathbb{E}_P[g(X)] - \mathbb{E}_Q\left[g(X) + \frac{g^2(x)}{4}\right].$

Solution. For Chi-squared distance we have that

$$\begin{aligned} f(x) &= (x-1)^2, \\ f^*(y) &= \sup_{x \in (0, \infty)} xy - (x-1)^2. \end{aligned}$$

Letting $s(x) = xy - (x-1)^2$ we compute the location of the maxima by

$$\begin{aligned} \left. \frac{ds(x)}{dx} \right|_{x>0} &= y - 2(x-1) \\ &\rightarrow x = \frac{y}{2} + 1 \end{aligned}$$

Thus

$$\begin{aligned} f^*(y) &= \frac{y^2}{2} + y - \left(\frac{y}{2} + 1 - 1\right)^2 \\ &= \frac{y^2}{4} + y \end{aligned}$$

and

$$\begin{aligned} \chi^2(P||Q) &= \sup_{g:\mathcal{X} \rightarrow \mathbb{R}} \mathbb{E}_P[g(x)] - \mathbb{E}_Q[f^*(g(x))] \\ &= \sup_{g:\mathcal{X} \rightarrow \mathbb{R}} \mathbb{E}_P[g(x)] - \mathbb{E}_Q\left[\frac{g^2(x)}{4} + g(x)\right] \end{aligned}$$

■

6. Inequalities Between f -Divergences.

Prove the following:

- (a) For any distributions $P, Q \in \mathcal{P}(\mathcal{X})$ it holds that

$$D_{\text{KL}}(P\|Q) \leq \log(1 + \chi^2(P\|Q)) \leq \chi^2(P\|Q).$$

Solution. Note that $\chi^2(P\|Q)$ is non-negative since it is an f -divergence. So we have $\chi^2(P\|Q) > -1$ and thus $\log(1 + \chi^2(P\|Q)) \leq \chi^2(P\|Q)$. Now we expand the log term:

$$\begin{aligned} \log(1 + \chi^2(P\|Q)) &= \log\left(1 + \int \frac{dQ}{d\lambda}(x) \left(\frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)} - 1\right) d\lambda\right) \\ &= \log\left(1 + \int \frac{dQ}{d\lambda}(x) \left(\frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)}\right)^2 d\lambda - 1\right) \\ &= \log\left(\int \frac{dP}{d\lambda}(x) \left(\frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)}\right) d\lambda\right) \end{aligned}$$

Now we expand and rewrite the KL divergence using the alternative form $D_{-\log(x)}(Q\|P)$:

$$\int \frac{dP}{d\lambda}(x) \left(-\log\left(\frac{\frac{dQ}{d\lambda}(x)}{\frac{dP}{d\lambda}(x)}\right)\right) d\lambda = \int \frac{dP}{d\lambda}(x) \left(\log\left(\frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)}\right)\right) d\lambda$$

Finally, we use Jansen's equality. Since \log is concave, we have $\mathbb{E}[\log(f(x))] \leq \mathbb{E}[\int f(x)]$. So we have

$$\begin{aligned} D_{\text{KL}}(P\|Q) &= \int \frac{dP}{d\lambda}(x) \left(\log\left(\frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)}\right)\right) d\lambda \\ &\leq \log\left(\int \frac{dP}{d\lambda}(x) \left(\frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)}\right) d\lambda\right) \\ &= \log(1 + \chi^2(P\|Q)) \end{aligned}$$

as desired. ■

- (b) Assume that $P = \text{Ber}(p)$ and $Q = \text{Ber}(q)$ where $p, q \in (0, 1)$. Show that

$$\delta_{\text{TV}}(P, Q)^2 \leq \frac{\ln(2)}{2} D_{\text{KL}}(P\|Q).$$

Solution. We expand both sides under the assumption that P and Q are Bernoulli measures. We get

$$f(p, q) = D_{\text{KL}}(P \| Q) - \frac{2}{\ln(2)} \delta_{\text{TV}}(P, Q)^2 = p \log\left(\frac{p}{q}\right) + (1-p) \log\left(\frac{1-p}{1-q}\right) - \frac{2}{\ln(2)} |p - q|^2$$

Taking the derivative, we get that the partials

$$\begin{aligned} \frac{\delta f}{\delta p} &= \frac{-\ln\left(\frac{p-1}{q-1}\right) + \ln\left(\frac{p}{q}\right) - 4p + 4q}{\ln(2)} \\ \frac{f}{\delta q} &= \frac{(1-2q)^2(p-q)}{(q-1)q \ln(2)} \end{aligned}$$

Note that the critical point is when $p = q$, which gives us zero for $f(p, p) = f(q, q) = 0$. For the inequality to hold, we need the critical point to be that of a local maxima. However, by taking the determinant of the hessian we would get that it's 0 so the second derivative test is inconclusive. So instead, we evaluate f by perturbing p , and have $p = q + \epsilon$, and $p = q - \epsilon$ and verify that $f(p, q) < f(q + \epsilon, q)$ and $f(p, q) < f(q - \epsilon, q)$. ■

(c) Assume that P and Q have finite supports. Show that

$$\delta_{\text{TV}}(P, Q)^w \leq \frac{1}{2} D_{\text{KL}}(P \| Q).$$

Solution. First note that the derivative of $(4+2x)h(x) - 3(x+1)^2$ is $-8x + 4(x+1) \log(x) + 8$, which is 0 at 1, and achieve a minimum value of 0.

observe that for $h(x) = x \log(x) + x - 1$, we have $\mathbb{E}_Q[h(\frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)})] = \mathbb{E}_Q[\frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)} \log(\frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)}) + \frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)} - 1] = \mathbb{E}_Q[\frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)} \log(\frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)})] = D_{\text{KL}}(P \| Q)$. Denote $X = \frac{\frac{dP}{d\lambda}(x)}{\frac{dQ}{d\lambda}(x)}$. Using the inequality in the hint, we have

$$\begin{aligned} D_{\text{KL}}(F) &\geq \frac{3}{2} \mathbb{E}_Q\left[\frac{(F-1)^2}{2+F}\right] \\ &= \frac{3}{2} \mathbb{E}_Q\left[\frac{(F-1)^2}{2+F}\right] \frac{1}{3} \mathbb{E}_Q[2+F] \quad \text{since } \frac{1}{3} \mathbb{E}_Q[2+F] = 1 \\ &\geq \frac{1}{2} \mathbb{E}_Q\left[\sqrt{\frac{(x-1)^2}{(2+x)(2+x)}}\right] \quad \text{using Cauchy-Schwarz} \\ &= \frac{1}{2} \mathbb{E}_Q[F-1]^2 \\ &= \frac{1}{2} \delta_{\text{TV}}(P, Q) \end{aligned}$$

■