# Statistical Distances

## Heuna Kim

# Examples of statistical distances

- Kolomogorov-Smirnov Statistic

- Kullback Leibler Divergence

- F-divergence and Bergmann divergence

- Jensen-Shannon Divergence

- Wasserstein Distance (Earthmover's distance)

- Mahalanobis Distance and Cook's distance

# Definition of Metrics

A metric on a set $X$ is a function (called *distance function* or simply *distance*)

$$d : X \times X \to [0, \infty),$$

where $[0, \infty)$ is the set of non-negative real numbers and for all $x, y, z \in X$, the following three axioms are satisfied
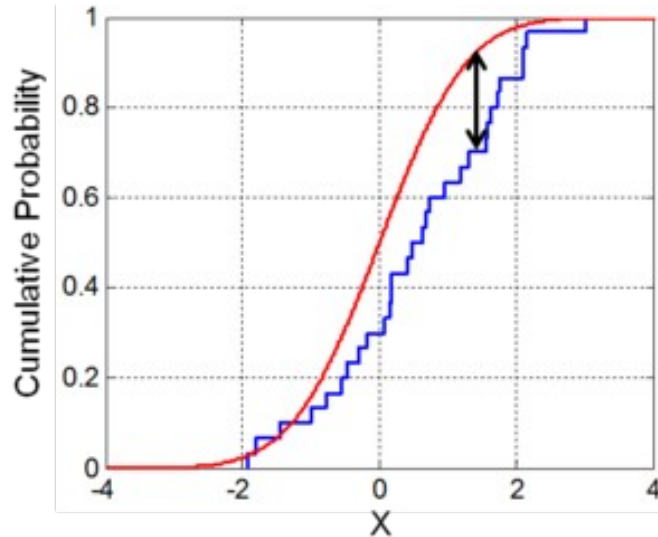
1. $d(x, y) = 0 \Leftrightarrow x = y$          identity of indiscernibles
2. $d(x, y) = d(y, x)$             symmetry
3. $d(x, y) \leq d(x, z) + d(z, y)$    triangle inequality

Statistical distances don't satisfy these all.

e.g.) cosine distance is a metric.

# Kolomogorov-Smirnov Statistic

- K-S Test: are two samples drawn from populations of the same distribution (univariate)?
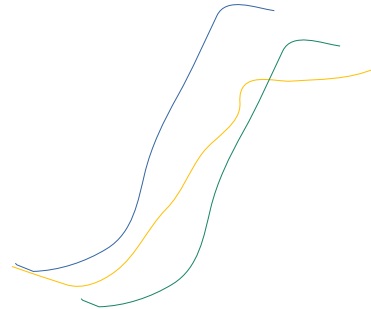


The Kolmogorov–Smirnov statistic for a given cumulative distribution function $F(x)$ is

$$D_n = \sup_x |F_n(x) - F(x)|$$

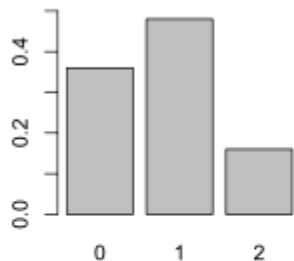$$\sqrt{n} D_n > K_\alpha,$$

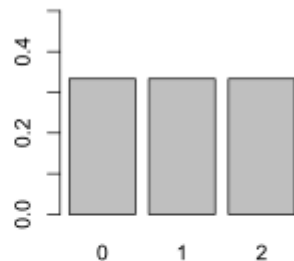where $K_\alpha$ is found from

$$\Pr(K \le K_\alpha) = 1 - \alpha.$$

# Total Variation Distance

$$\delta(P, Q) = \sup_{A \in \mathcal{F}} |P(A) - Q(A)|.$$

$\mathcal{F}$: all possible events (subsets of a sample space)

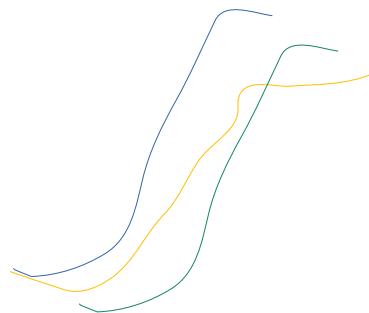**Distribution P**
**Binomial with p = 0.4 , N = 2**

**Distribution Q**
**Uniform with p = 1/3**

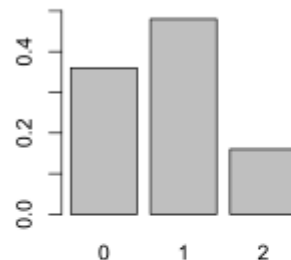| x | 0 | 1 | 2 |
|---|---|---|---|
| Distribution $P(x)$ | $9/25$ | $12/25$ | $4/25$ |
| Distribution $Q(x)$ | $1/3$ | $1/3$ | $1/3$ |

# Kullback Leibler Divergence

- Relative entropy or information gain
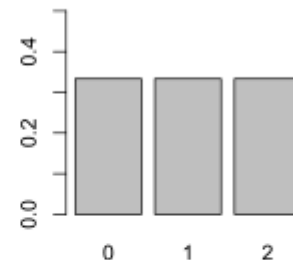
$$D_{\mathrm{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log\left(\frac{P(x)}{Q(x)}\right).$$

$$= -\sum_{x \in \mathcal{X}} p(x) \log q(x) + \sum_{x \in \mathcal{X}} p(x) \log p(x)$$

$$= \mathrm{H}(P, Q) - \mathrm{H}(P)$$

$$D_{\mathrm{KL}}(P \parallel Q) = \int_{-\infty}^{\infty} p(x) \log\left(\frac{p(x)}{q(x)}\right) dx$$



**Distribution P**
Binomial with p = 0.4 , N = 2

**Distribution Q**
Uniform with p = 1/3

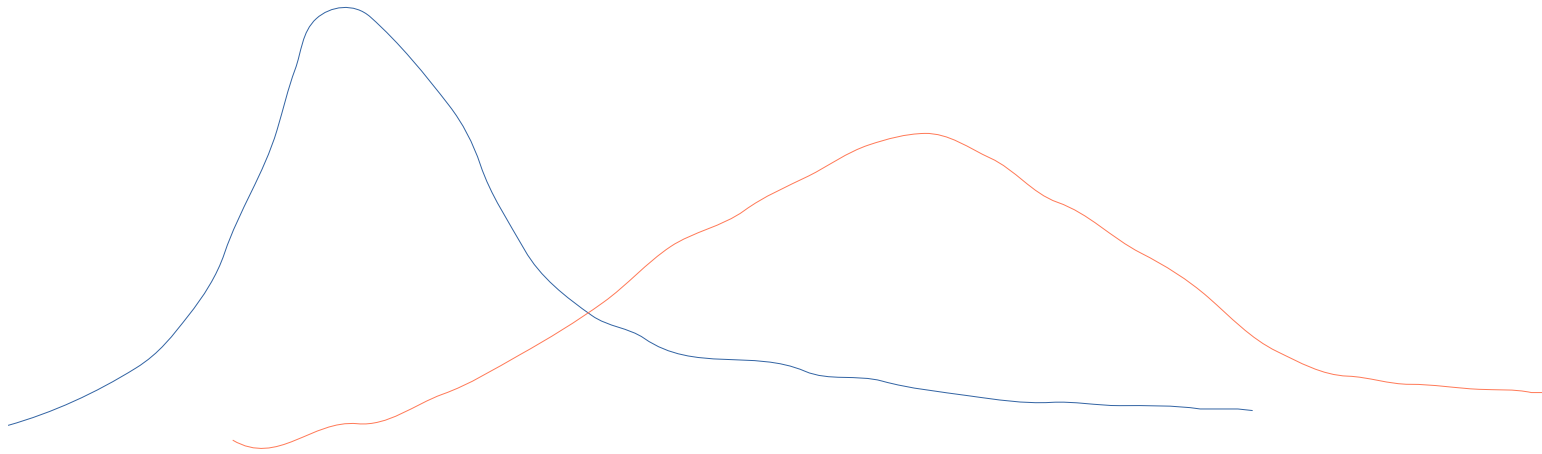| x | 0 | 1 | 2 |
|---|---|---|---|
| Distribution $P$(x) | 9/25 | 12/25 | 4/25 |
| Distribution $Q$(x) | 1/3 | 1/3 | 1/3 |

# To Note about KL-Divergence

- Not a metric
  – Why?
- f-divergence $\longrightarrow$ $D_f(P \parallel Q) \equiv \int_\Omega f\left(\dfrac{dP}{dQ}\right) dQ.$
- Bergman-divergence

  • **Convexity:** $D_F(p, q)$ is convex in its first argument, but not necessarily in the second argument

# Exercise

# Jensen-Shannon divergence

- A symmetrized version of KL divergence

$$D_{\mathrm{JS}} = \frac{1}{2} D_{\mathrm{KL}}\left(P \parallel M\right) + \frac{1}{2} D_{\mathrm{KL}}\left(Q \parallel M\right)$$

where $M$ is the average of the two distributions,

$$M = \frac{1}{2}(P + Q).$$

# Adversarial loss and JS div

$$\min_{G}\max_{D} L(D, G) = \mathbb{E}_{x \sim p_r(x)}[\log D(x)] + \mathbb{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))]$$

$$= \mathbb{E}_{x \sim p_r(x)}[\log D(x)] + \mathbb{E}_{x \sim p_g(x)}[\log(1 - D(x)]$$

When both $G$ and $D$ are at their optimal values, we have $p_g = p_r$ and $D^*(x) = 1/2$ and the loss function becomes:

$$L(G, D^*) = \int_x \left( p_r(x) \log(D^*(x)) + p_g(x) \log(1 - D^*(x)) \right) dx$$

$$= \log \frac{1}{2} \int_x p_r(x) dx + \log \frac{1}{2} \int_x p_g(x) dx$$

$$= -2 \log 2$$

$$D_{JS}(p_r \| p_g) = \frac{1}{2} D_{KL}(p_r \| \frac{p_r + p_g}{2}) + \frac{1}{2} D_{KL}(p_g \| \frac{p_r + p_g}{2})$$

$$= \frac{1}{2} \left( \log 2 + \int_x p_r(x) \log \frac{p_r(x)}{p_r + p_g(x)} dx \right) +$$

$$\frac{1}{2} \left( \log 2 + \int_x p_g(x) \log \frac{p_g(x)}{p_r + p_g(x)} dx \right)$$

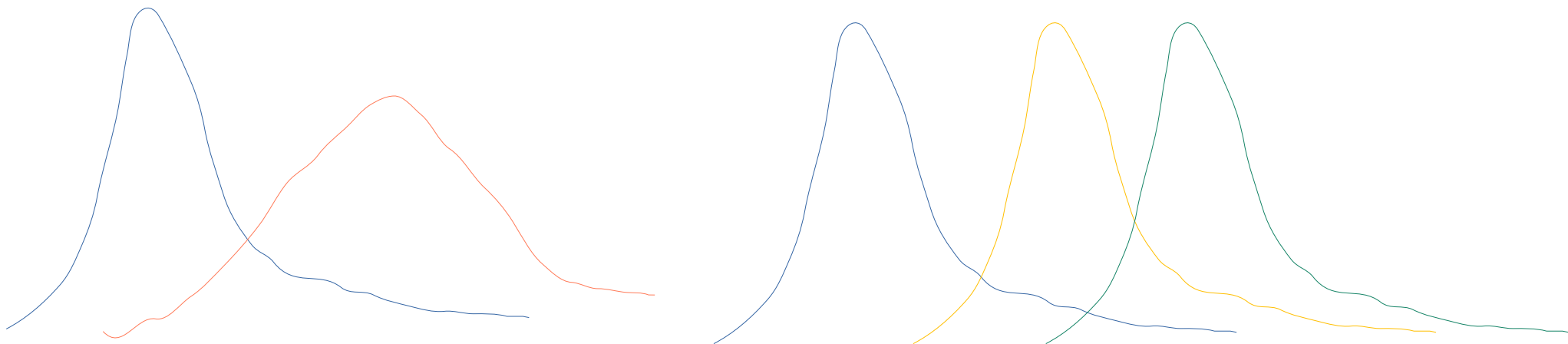$$= \frac{1}{2} \left( \log 4 + L(G, D^*) \right)$$

# Wasserstein Metric

- Earthmover's distance and optimal transport plan

The $p^{\text{th}}$ **Wasserstein distance** between two probability measures $\mu$ and $\nu$ in $P_p(M)$ is defined as

$$W_p(\mu, \nu) := \left( \inf_{\gamma \in \Gamma(\mu, \nu)} \int_{M \times M} d(x, y)^p \, \mathrm{d}\gamma(x, y) \right)^{1/p},$$

where $\Gamma(\mu, \nu)$ denotes the collection of all measures on $M \times M$ with marginals $\mu$ and $\nu$ on the first and second factors respectively.

# Kantorovich-Rubinstein duality

If $\|f\|_L \leq K,$ and both distributions have bounded supports

$$W(p_r, p_g) = \frac{1}{K} \sup_{\|f\|_L \leq K} \mathbb{E}_{x \sim p_r}[f(x)] - \mathbb{E}_{x \sim p_g}[f(x)]$$

# To note about Wasserstein distance

- Metric

- If two distribtions are multidimentional gaussian,
  equivalent to frechet inception distance.

$$\text{FID} = |\mu - \mu_w|^2 + \text{tr}(\Sigma + \Sigma_w - 2(\Sigma\Sigma_w)^{1/2}).$$

- Cf. frechet distance == dogwalker's distance
  - For similarity of two curves

# Mahalanobis distance

- Unit-less scale-invariant metric

$$d(\vec{x}, \vec{y}) = \sqrt{(\vec{x} - \vec{y})^{\top} \mathbf{S}^{-1} (\vec{x} - \vec{y})}.$$

- Bergmann divergence

- Assume bell shapes

- cf. Cook's distance:

$$D_i = \frac{\sum_{j=1}^{n} \left( \widehat{y}_j - \widehat{y}_{j(i)} \right)^2}{ps^2}$$

$$\underset{n \times 1}{\mathbf{y}} = \underset{n \times p}{\mathbf{X}} \ \underset{p \times 1}{\boldsymbol{\beta}} \ + \ \underset{n \times 1}{\boldsymbol{\varepsilon}}$$

where $\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \sigma^2 \mathbf{I}\right)$ is the error term, $\boldsymbol{\beta} = [\beta_0 \, \beta_1 \ldots \beta_{p-1}]$

where $\widehat{y}_{j(i)}$ is the fitted response value obtained when excluding $i$, and $s^2 = \dfrac{\mathbf{e}^{\top}\mathbf{e}}{n-p}$