

KL Divergence,
Fano

- Kullback-Leibler divergence
- Fano's inequality for classifiers

Big Picture

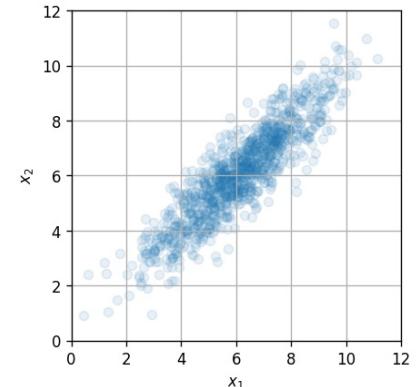
- *unsupervised* machine learning is about finding interesting patterns in data

$$\mathcal{D} = \{(x_i)\}_{i=1}^n$$

- learn $p(\mathbf{x})$

- how well can we approximate $p(\mathbf{x})$ or $p(\mathbf{x}, y)$?

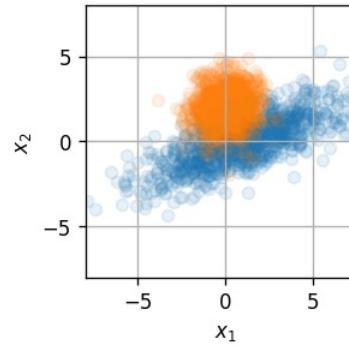
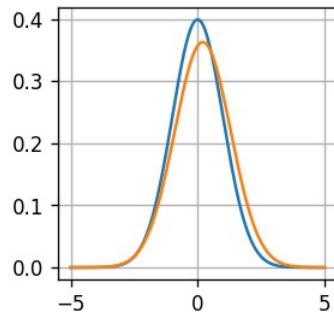
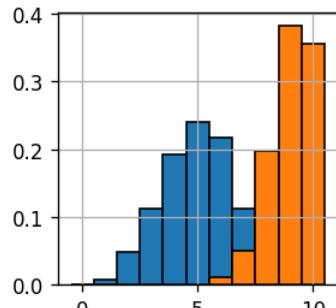
- simpler question - how well can we estimate $E[\mathbf{x}]$?



- intuition: more samples, better job learning $p(\mathbf{x})$.

Kullback-Leibler Divergence

- basic task - measuring how different/similar two distributions are:



- many choices: ℓ_1 , ℓ_2 , KL divergence, total variation, Bhattacharyya distance, Hellinger, ...
- the *Kullback-Leibler Divergence* between two distributions p and q is:

If $p(x) > 0$ implies $q(x) > 0$, then

$$D(p||q) = \sum p(x) \log \frac{p(x)}{q(x)}$$

$$D(p||q) = \int p(x) \log \frac{p(x)}{q(x)} dx$$

where the sum or integral is over x with $p(x) > 0$. Otherwise, $D(p||q) = \infty$.

Example

- example: compute $D(p||q)$ when $p \sim \mathcal{N}(0, 1)$ and $q \sim \mathcal{N}(\mu, 1)$

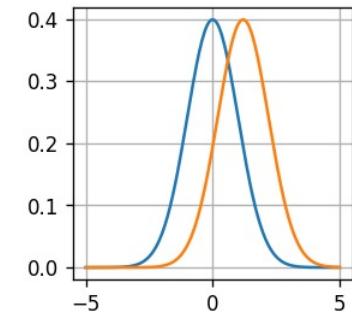
$$D(p||q) = E_p \left[\log \frac{p(x)}{q(x)} \right]$$

$$D(p||q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx$$

$$= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \log \frac{e^{-x^2/2}}{e^{-(x-\mu)^2/2}} dx$$

$$= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \left(\frac{\mu^2}{2} - x\mu \right) dx$$

$$= \frac{\mu^2}{2} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx + \mu \int_{-\infty}^{\infty} x \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = \frac{\mu^2}{2}$$



Kullback-Leibler Divergence

- the *Kullback-Leibler Divergence* between two distributions p and q is:

$$D(p||q) = E_p \left[\log \frac{p(x)}{q(x)} \right]$$

- why is $D(p||q)$ important?

- imagine a classification task between $X_i \stackrel{i.i.d.}{\sim} p$ vs. $X_i \stackrel{i.i.d.}{\sim} q$, and the maximum likelihood rule:

$$\log \frac{p(x)}{q(x)} \leq 0 \quad \sum_{i=1}^n \log \frac{p(x_i)}{q(x_i)} \leq 0$$

- $D(p||q)$ is the expected value of the log likelihood ratio (when p is true class)

Chernoff-Stein Lemma:

Consider a simple hypothesis in which you observe n i.i.d. samples from either p or q , and any decision rule that has $\mathbb{P}(\text{decide } q|p \text{ true}) \leq \epsilon$, then

$$\mathbb{P}(\text{decide } p|q \text{ true}) \gtrsim e^{-nD(p||q)}$$

- fix one type of error, and the other error decays exponentially in n and $D(p||q)$.

Kullback-Leibler Divergence

- the *Kullback-Leibler Divergence* between two distributions p and q is:

$$D(p||q) = \sum p(x) \log \frac{p(x)}{q(x)}$$

- why is $D(p||q)$ important?

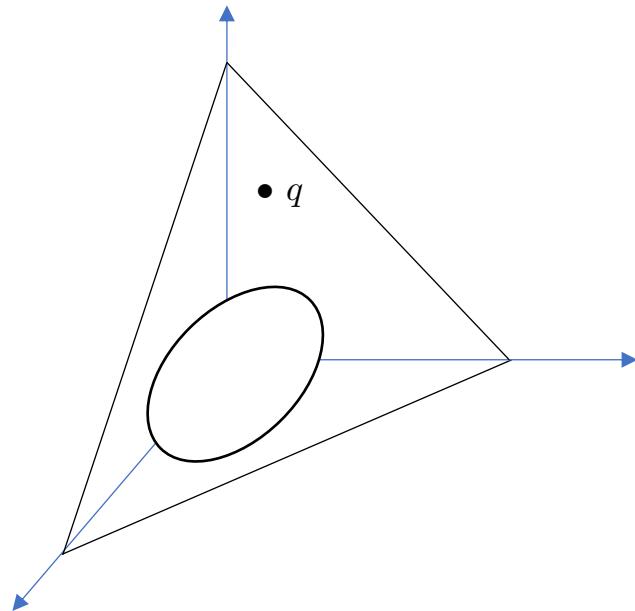
- can be used to find distribution *closest* to prior assumption and observations

- 1) assume a prior distribution q
- 2) make observations that constrain problem

$$\min_p D(p||q)$$

subject to p agrees with constraints

- $D(p||q)$ is convex in both p and q
- we'll do this when we look at variational autoencoders



Properties of KL Divergence

- the *Kullback-Leibler Divergence* between two distributions p and q is:

$$D(p||q) = E_p \left[\log \frac{p(x)}{q(x)} \right]$$

- KL divergence is non-negative: $D(p||q) \geq 0$ and $D(p||q) = 0$ if and only if $p(x) = q(x)$ for all x with $p(x) > 0$

sketch of proof:

$$\begin{aligned} -D(p||q) &= - \sum p(x) \log \left(\frac{p(x)}{q(x)} \right) \\ &= \sum p(x) \log \left(\frac{q(x)}{p(x)} \right) \\ &\leq \log \left(\sum p(x) \frac{q(x)}{p(x)} \right) \quad \text{Jensen: } E[f(X)] \leq f(E[X]) \text{ for } f(x) \text{ concave} \\ &= 0 \end{aligned}$$

Jensen's Inequality

Jensen's inequality: for a convex function $f(\cdot)$ and a random variable X

$$E[f(X)] \geq f(E[X])$$

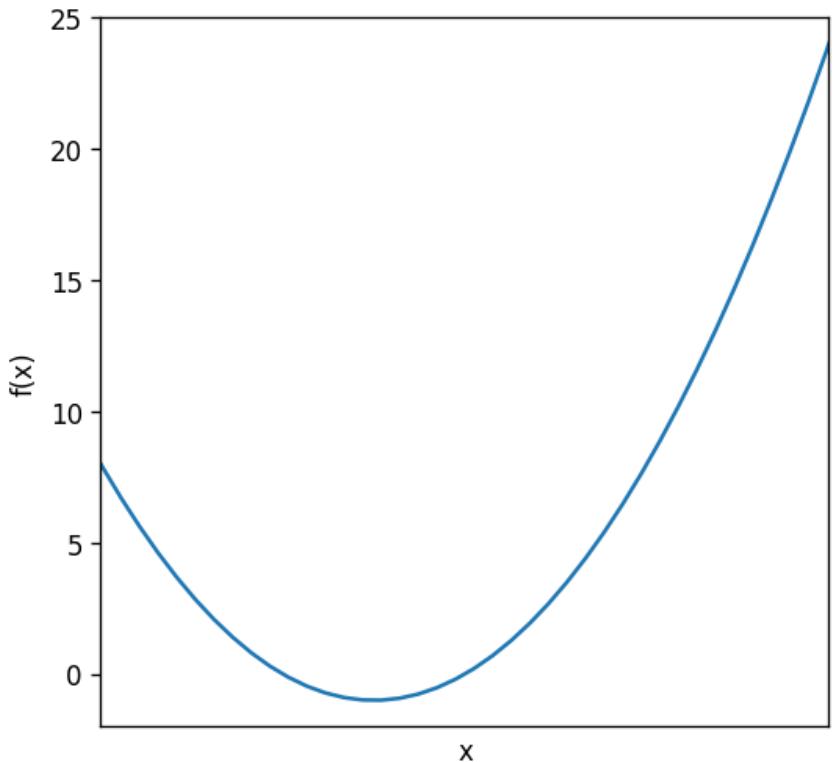
$$\sum p(x)f(x) \geq f(\sum xp(x))$$

sketch of proof:

1) consider binary X

$$p(x_1)f(x_1) + p(x_2)f(x_2) \geq f(x_1p(x_1) + x_2p(x_2))$$

2) by induction



Alternative Derivation of Jensen's Inequality

Recall that f is convex if for $0 \leq \lambda \leq 1$,

$$f(u + \lambda(v - u)) \leq f(u) + \lambda[f(v) - f(u)].$$

For $0 < \lambda \leq 1$, this is equivalent to

$$\frac{f(u + \lambda(v - u)) - f(u)}{\lambda} \leq f(v) - f(u).$$

Letting $\lambda \rightarrow 0$ implies

$$\nabla f(u)^T(v - u) \leq f(v) - f(u),$$

or

$$f(u) + \nabla f(u)^T(v - u) \leq f(v).$$

Now suppose X is a random vector with mean μ . Put $u = \mu$ and $v = X$ to get

$$f(\mu) + \nabla f(\mu)^T(X - \mu) \leq f(X).$$

Taking expectations yields $f(\mu) \leq \mathbb{E}[f(X)]$, or

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)].$$

Entropy, Mutual Information, and KL Divergence

- definition - entropy:

$$H(X) = E \left[\log_2 \left(\frac{1}{p(x)} \right) \right]$$

- relating conditional entropy to joint (chain rule):

$$H(X, Y) = H(X) + H(Y|X)$$

- amount of uncertainty in a random variable

- mutual information: reduction in uncertainty of Y due to knowledge of X

$$I(Y; X) = H(Y) - H(Y|X)$$

$$I(Y; X) = E \left[\log \left(\frac{p(x, y)}{p(y)p(x)} \right) \right]$$

- KL divergence:

$$D(p||q) = E_p \left[\log \frac{p(x)}{q(x)} \right]$$

- mutual information is equal to the KL divergence between the joint and the product of marginals

$$I(X; Y) = D(p(x, y)||p(x)p(y))$$

Properties of Entropy

- definition - entropy:

$$H(X) = E \left[\log_2 \left(\frac{1}{p(x)} \right) \right]$$

- properties of entropy:

$$H(X) \geq 0$$

$$H(X) \leq \log_2(|\mathcal{X}|)$$

proof:

$$D(p||q) = \sum p(x) \log \frac{p(x)}{q(x)} \quad \text{set } q(x) = \frac{1}{|\mathcal{X}|}, \text{ let } X \sim p(X)$$

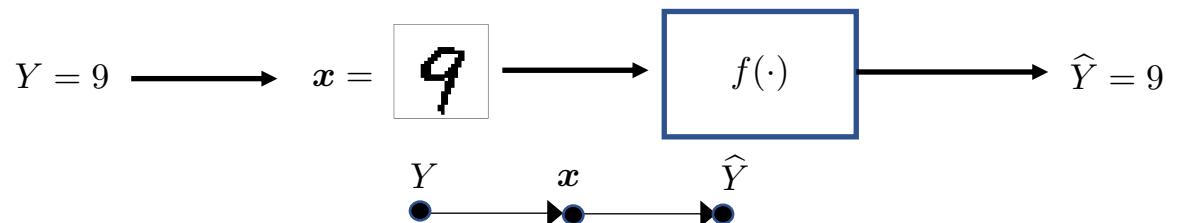
$$= \sum p(x)(\log |\mathcal{X}| + \log p(x))$$

$$= \log |\mathcal{X}| - H(X)$$

$$\geq 0$$

$$H(X) \leq \log |\mathcal{X}| \text{ with equality when } p(x) = \frac{1}{|\mathcal{X}|}$$

Fano's Inequality



- for any $\hat{Y} = f(x)$ we have $Y \rightarrow x \rightarrow \hat{Y}$ and

$$\mathbb{P}(Y \neq \hat{Y}) \geq \frac{H(Y|x) - 1}{\log(|\mathcal{Y}|)}$$

where $|\mathcal{Y}|$ is the number of classes

- proof: define $E = \mathbb{I}_{\{\hat{Y} \neq Y\}}$

$$\begin{aligned} H(E, Y|\hat{Y}) &= H(Y|\hat{Y}) + H(E|Y, \hat{Y}) \\ &= H(Y|\hat{Y}) \\ &\geq H(Y|x) \end{aligned}$$

$$\begin{aligned} H(E, Y|\hat{Y}) &= H(E|\hat{Y}) + H(Y|E, \hat{Y}) \\ &\leq 1 + H(Y|E, \hat{Y}) \\ &\leq 1 + \mathbb{P}(E=1)H(Y|E=1, \hat{Y}) + \mathbb{P}(E=0)H(Y|E=0, \hat{Y}) \\ &\leq 1 + \mathbb{P}(\hat{Y} \neq Y)H(Y|E=1, \hat{Y}) \\ &\leq 1 + \mathbb{P}(\hat{Y} \neq Y)\log|\mathcal{Y}| \end{aligned}$$