

Lecture 5 : 19th Jan 2026 :

Distributional Divergence

Objective: Given a pair of distribution func defined
on the same sample space, define / quantify
"distance" b/w them.

Given an event $A \in \mathcal{F}$,

$$\text{define } I(A) \triangleq -\log IP(A)$$

Likely info has less info

less likely info has more info

→ Info of Ω has zero info

→ null event has ∞ info

→ A & B are indpt, $IP(AB) = IP[A] IP[B]$

$$\log IP[AB] = \log IP[A] + \log IP[B]$$

→ $-\log IP(A)$ satisfies all the above Cond'

→ To transfer the info, how many bits I need.

* Suppose X is DRV

then the "Surprisal" associated with it is

quantified as follows

$$\begin{aligned} & \mathbb{E}_{P_x} \left[-\log p_x(x) \right] \quad \text{by LOTUS.} \\ &= \sum_{i=1}^n p_x(x_i) \log p_x(x_i) \quad : \text{Avg info in } P_x \\ & \qquad \qquad \qquad \text{Expected Value of } P_x \end{aligned}$$

* This is the entropy of P_x

$$H(P_x) = \mathbb{E}_{P_x} \left[-\log p_x(x) \right]$$

Suppose P_x & Q_x are 2 distributions on the

same sample space

$$\text{Define } \mathbb{E}_{P_x} -\log Q_x(x) = H(P_x, Q_x) : \text{Cross entropy}$$

Samples from P_x , estimate Q_x

Avg info Q_x contain about P_x

* info that one distribution contains about another distribution.

* Cross entropy is not symmetric.

$$\text{Consider } -H(P_x) + H(P_x, Q_x)$$

$$= + \sum p_x(x_i) \log p_x(x_i) - \sum p_x(x_i) \log q_x(x_i)$$

Extra info of Q_x which is present

$$D_{KL}(P_x \parallel Q_x) = \sum p_x(x_i) \log \frac{p_x(x_i)}{q_x(x_i)}$$

Show:

$$\Rightarrow D_{KL}(P_x \parallel Q_x) \neq D_{KL}(Q_x \parallel P_x)$$

$$2) D_{KL} \geq 0$$

$$3) D_{KL} = 0 \text{ if } P_x = Q_x$$

Does not obey Law of triangularity so its not a metric.

Now a similar defn when X is CRV

$$D_{KL}(P_x \parallel Q_x) = \int_X p_x(x) \log \frac{p_x(x)}{q_x(x)} dx.$$

Suppose we have "N" Samples drawn iid.

$$D = \{v_1, v_2, \dots, v_n\} \sim \text{iid } P_v$$

Goal is to estimate P_v given D.

Assume $f_\theta(v)$ as the model density & compute

$$D_{KL}(P_v \parallel P_\theta) = \int f_v(v) \log \frac{f_v(v)}{f_\theta(v)} dv.$$

$$\theta^* = \underset{\theta}{\operatorname{argmin}} D_{KL}(P_v \parallel P_\theta)$$

$$= \underset{\theta}{\operatorname{argmin}} \left[\int f_v(v) \log \frac{f_v(v)}{f_\theta(v)} dv - \int f_v(v) \log f_\theta(v) dv \right] \xrightarrow{\text{Independent of } \theta}$$

$$= \underset{\theta}{\operatorname{argmin}} \left[- \int f_v(v) \log f_\theta(v) dv \right]$$

$$= \underset{\theta}{\operatorname{argmin}} - \int f_v(v) \log f_\theta(v) dv$$

$$= \underset{\theta}{\operatorname{argmin}} - \mathbb{E}_{P_v} \log f_\theta(v)$$

wkt expectations can be approximated via samples

By the Law of Large numbers.

$$\underset{IP_V}{\mathbb{E}} \log p_{\theta}(v) \approx \frac{1}{N} \sum_{i=1}^n \log p_{\theta}(v_i)$$

where $v_i \sim i.i.d \text{ } IP_V$

$$\frac{1}{N} \sum_{i=1}^N \log p_{\theta}(v_i) \xrightarrow{\text{IP}_V} \underset{IP_V}{\mathbb{E}} \log p_{\theta}(v) \text{ as } N \rightarrow \infty$$

Now our optimization problem is,

Negative
Log Likelihood
(NLL)

$$\therefore \theta^* = \underset{\theta}{\operatorname{argmin}} \left[\frac{1}{N} \sum_{i=1}^N \log p_{\theta}(v_i) \right]$$

This is the Minimal DKL estimator.

or we say it as Maximum Likelihood
estimator

Another Algebra : Joint Likelihood of the Data

$$L(\theta) = \prod_{i=1}^n p_{\theta}(v_i)$$

taking a monotonic log func

$$L_\theta(D) = \sum_{i=1}^N \log p_\theta(\omega_i)$$

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \frac{1}{N} \sum_{i=1}^N \log p_\theta(\omega_i)$$

How to solve the ML estimation problem in practice?

a) how to choose p_θ ?

b) how to solve the optimization problem?

Solution to the problems:

a) p_θ is chosen using a Canonical func approx

e.g.: Assume $p_\theta(\cdot) \sim N(\cdot; h_\theta(v), I)$

$h_\theta(v)$: hypothesis func/Model

→ Now it boils to estimating parameters of this func

b) Optimization is mostly numerical

Consider the problem of supervised regression.

$$x \in \mathbb{R}^d, y \in \mathbb{R}$$

Data $D = \{(x_i, y_i)\}_{i=1}^N \sim \text{iid } P_{xy}$

Estimate $P_{y|x}$

$$p_\theta(y|x) \sim N(y; h_\theta(x), I)$$

Assumption.
Model

$$h_\theta(x) : x \rightarrow \mathbb{R}$$

$$\theta^* = \underset{\theta}{\operatorname{argmin}} -\frac{1}{N} \sum_{i=1}^N \log p_\theta(y_i|x_i)$$

$$\propto \underset{\theta}{\operatorname{argmin}} -\frac{1}{N} \sum_{i=1}^N \log \exp \{- (y_i - h_\theta(x_i))^2\}$$

$$\theta^* \propto \underset{\theta}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N (y_i - h_\theta(x_i))^2$$

This is the least square estimate

Lecture 6: 21st Jan 2026:

Risk Minimization Framework

Given Data $D = \{(x_i, y_i)\}_{i=1}^N \sim \text{iid } \mathbb{P}_{X,Y}$
 $x \in \mathbb{R}^d, y \in \mathbb{R}^k / \{1, \dots, k\}$

Goal is to estimate hypothesis func

$h : X \rightarrow Y$ $h \in \mathcal{H}$ set of all func
 from $X \rightarrow Y$.

Quantify the "goodness" of an h

Define a notion of Loss-func L

$$L : Y \times Y \rightarrow \mathbb{R}^+$$

$$L(h(x), y) \rightarrow \mathbb{R}^+$$

$$\text{Eg: a) } L(h(x), y) = \|h(x) - y\|_2^2$$

squared error Loss.

$$\text{b) } L(h(x), y) = \begin{cases} 0 & \text{if } h(x) = y \\ 1 & \text{if } h(x) \neq y \end{cases}$$

Zero-One Loss.

Extending the Loss func for the population.

Define $R(h) \triangleq \mathbb{E}_{P_{xy}} L(h(x), y)$ Risk function.

Problem : Estimate $h^*(x)$

$$h^*(x) = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \mathbb{E}_{P_{xy}} L(h(x), y)$$

This is known as the Risk minimization framework.

We can't solve this as we don't know P_{xy}

∴ we look @ its surrogate

True risk can't be minimized ∵ P_{xy} is Unknown

∴ Consider the Surrogate

$$\hat{R}(h) \triangleq \frac{1}{n} \sum_{i=1}^n L(h(x_i), y_i) \quad : \text{Empirical Risk,}$$

$$\hat{h}^*(x) = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \hat{R}(h)$$

$x_i, y_i \sim \text{iid } P_{xy}$

This is Empirical Risk minimization (ERM_s)

We know

$$\hat{R}(h) \rightarrow R(h) : \text{Law of Large numbers.}$$

However $\hat{h}^*(x)$ may not converge $h^*(x)$

Central question : when would $\hat{h}^*(x) \rightarrow h^*(x)$?

Example of ERM for a Regression problem:

Given $D = \{(x_i, y_i)\}_{i=1}^N \sim \text{iid } P_{xy}$

$x_i \in \mathbb{R}^d, y \in \mathbb{R}$

Suppose $h_\theta(x)$ comes from a family of parametric func parameterized by θ

$$h_\theta(x) = \theta^T x @ h_\theta(x) = e^{\theta^T x}$$

Define Loss func mean sq error

$$\hat{R}(h_\theta) = \frac{1}{N} \sum_{i=1}^N (h_\theta(x_i) - y_i)^2$$

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \hat{R}(h_\theta)$$

The above is equivalent to estimating $f_\theta(y|x)$

using ML estimation with

$$f_\theta(y|x) \sim N(y; h(x), \Sigma)$$

ERM @ MLE is Solving same problem

Transformers : Solve ERM

GMM : Solve MLE

Diffusion model

We showed the equivalence in optimization way now

Let's look at its equivalence in solution space.

Equivalence between ERM & Divergence Minimization (MLE)

Question: what func minimizes the true risk?

Consider $D = \{(x_i, y_i)\}_{i=1}^N \sim \text{iid } P_{xy}$

$x_i \in \mathbb{R}^d, y \in \{-1, 1\}$

$$L(h(x), y) = \begin{cases} 0 & \text{if } h(x) = y \\ 1 & \text{otherwise} \end{cases}$$

$$R(h) = \mathbb{E}_{P_{xy}} [L(h(x), y)]$$

Question: what h , minimizes the true risk for the above scenario?

Ans: Consider a hypothesis func as follows.

$$h_B(x=x) = \begin{cases} 0 & \text{if } p(y=0|x) > p(y=1|x) \\ 1 & \text{if } p(y=1|x) \geq p(y=0|x) \end{cases}$$

Bayes' Classifier

$$\text{Claim } R(h_B) \leq R(h) \quad \forall h \in \mathcal{H}$$

$$\text{Proof: Let } S_i(h) = \{x \in \mathbb{R}^d : h(x) = i\} \quad i = 0, 1$$

$$S_0(h) \cap S_1(h) = \emptyset$$

$$S_0(h) \cup S_1(h) = \mathbb{R}^d \quad \forall h \in \mathcal{H}$$

$$\text{Note: } R(h) = \mathbb{P}(I_{h(x) \neq y})$$

$$R(h) = \mathbb{E}_{P_{xy}} [I_{h(x) \neq y}]$$

$$= \mathbb{P}[I_{h(x) \neq y}]$$

$$\begin{aligned}
 &= \text{IP}[h(x)=1, Y=0] + \text{IP}[h(x)=0, Y=1] \\
 &= \text{IP}[x \in S_1(h), Y=0] + \text{IP}[x \in S_0(h), Y=1] \\
 &= \text{IP}[Y=0] \text{IP}_{x|y=0} [x \in S_1(h) | Y=0] + \\
 &\quad \text{IP}[Y=1] \text{IP}_{x|y=1} [x \in S_0(h) | Y=1]
 \end{aligned}$$

$$\begin{aligned}
 &= \text{IP}[Y=0] \int_{S_1(h)} p_{x|y=0}(x | y=0) dx + \\
 &\quad \text{IP}[Y=1] \int_{S_0(h)} p_{x|y=1}(x | y=1) dx
 \end{aligned}$$

→ for each example Only One integral will be evaluated. Since $S_1(h) \cap S_0(h) = \emptyset$.

→ ∀ $x \in \mathbb{R}^d$ Only one of two will be evaluated.

Trill now its true for every hypothesis.

Now for the specific case of $h_B(x)$
we have

$$R(h) = \int_{S_1} p_{y=0} p_{x|y=0} dx + \int_{S_0} p_{y=1} p_{x|y=1} dx$$

Recall

$$h_B = \begin{cases} 1 & \text{if } P_{y=1|x} > P_{y=0|x} \\ 0 & \text{otherwise} \end{cases}$$
$$\Rightarrow P_{x|y=1} P_{y=1} > P_{x|y=0} P_{y=0}$$

for the Bayes classifier the one with the minimal value will be evaluated.

$$R(h_B) = \int_{\mathbb{R}^D} \min(P_{y=0} P_{x|y=0}, P_{y=1} P_{x|y=1}) dx$$

$$\therefore R(h_B) \leq R(h) \quad \forall h \in \mathcal{H}$$

You can't design a classifier which is better than Bayes classifier.

(k-class classification
Some result)

Regression

Loss: S_2 error loss

Optimal hypothesis is $\mathbb{E}_{P_{y|x}}$

