

B. Discrete Valued RV. PMF

$\underline{X} \in \{\underline{x}_1, \dots, \underline{x}_c\} \sim P(\underline{x})$: True PMF
 $Q(\underline{x})$: Approx for $P(\underline{x})$

$$D_{KL}(P||Q) = E_{\underline{X} \sim P} \left[\ln \frac{P(\underline{X})}{Q(\underline{X})} \right]$$
$$= \sum_{i=1}^N P(x_i) \ln \frac{P(x_i)}{Q(x_i)}$$

Properties of KLD:

P1. Non-negative : $D_{KL}(P||Q) \geq 0 \forall P, Q$

P2: Equality:

$$D_{KL}(P||Q) = 0 \text{ iff } P(\underline{x}) = Q(\underline{x}) \forall \underline{x}$$

Proof of "sufficient": $\ln \frac{P(\underline{x})}{Q(\underline{x})} = \ln 1 = 0 \forall \underline{x}$

P1+P2:

KLD is a suitable metric for approximating P by Q .

P3. Asymmetric:

$$D_{KL}(P||Q) \neq D_{KL}(Q||P)$$
$$= E_{\underline{X} \sim P} \ln \left(\frac{P(\underline{x})}{Q(\underline{x})} \right) \neq E_{\underline{X} \sim Q} \ln \left(\frac{Q(\underline{x})}{P(\underline{x})} \right)$$

Hence it is a divergence, not distance

forward KLD

backward KLD

D_{KL} is NOT a true distance measure, with

$$D(\underline{x}, \underline{y}) \neq D(\underline{y}, \underline{x}).$$


We mainly use forward D_{KL} .


To minimize D_{KL} , always increase the denominator in the $\ln(\cdot)$ term.

$$\left\{ \ln \left(\frac{p(x)}{q(x)} \right) \right\} \downarrow$$

Exercise

• KLD between Gaussian and Laplace distribution

$$p(x) \sim N(0, \sigma^2), \quad p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}x^2}$$


$$q(x) \sim \text{Laplace}(0, b), \quad q(x) = \frac{1}{2b} e^{-\frac{|x|}{b}}$$


We can't get $D_{KL} = 0$ as both are different

Choose b such as to best approximate $p(x)$ by $q(x)$.

$$\frac{p(x)}{q(x)} = \frac{2b}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2} + \frac{|x|}{b}\right)$$

$$\ln\left(\frac{p(x)}{q(x)}\right) = \ln\left(\sqrt{\frac{2}{\pi}} \frac{b}{\sigma}\right) - \frac{x^2}{2\sigma^2} + \frac{|x|}{b}$$

$$D_{KL}(p||q) = E_{x \sim p} \ln\left(\frac{p(x)}{q(x)}\right) = \ln\left(\sqrt{\frac{2}{\pi}} \frac{b}{\sigma}\right) + E_{x \sim p} \left(-\frac{x^2}{2\sigma^2} + \frac{|x|}{b}\right)$$

$$E_{x \sim p}(x^2) = \sigma^2$$

$$E_{x \sim p}(|x|) = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2}\right) \cdot |x| dx$$

$$= \sqrt{\frac{2}{\pi}} \frac{\sigma}{b}$$

$$\text{If } \alpha = \frac{\sigma}{b}, \text{ Then } D_{KL}(p||q) = \ln\left(\sqrt{\frac{2}{\pi}}\right) - \ln \alpha - \frac{1}{2} + \sqrt{\frac{2}{\pi}} \alpha$$

Now we need to min α .

$$\frac{\partial D_{KL}}{\partial \alpha} = -\frac{1}{\alpha} + \sqrt{\frac{2}{\pi}} \stackrel{!}{=} 0$$

$$\Rightarrow \alpha = \sqrt{\frac{\pi}{2}}$$

$$\approx b \approx 0.8 \sigma$$

$$D_{KL}(p||q) \Big|_{\alpha=\sqrt{\frac{\pi}{2}}} = -\ln \frac{\pi}{2} + \frac{1}{2} \approx 0.048$$

Deep Learning

①

Assignments - Online - Mandatory for Lab.

"A recipe ... " - VVS - download and read.

1. Introduction

Slide 1.0

1.1 What is machine learning?

Slide 1.2-1.3

Signal to signal \rightarrow noised audio to denoised audio

" to parameter \rightarrow auto driving eg range of radar

" to class \rightarrow classification \rightarrow finite no of categories

} Regression
(continuous)

see radar workings papers \rightarrow from signal $x(n)$ to distⁿ of target
Sync using correlation.

Slide 1.9 - Overfitting is to be avoided so divide training ds to
Training + validation.

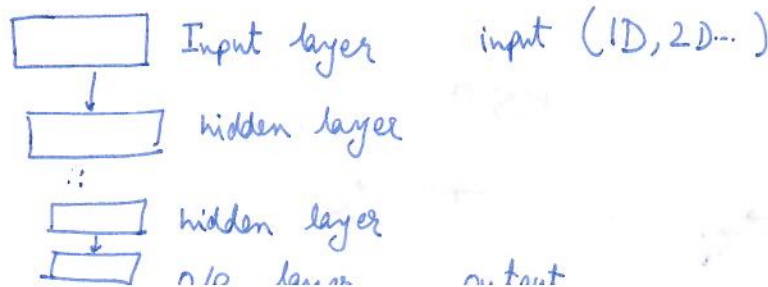
1.10 \rightarrow mainly we will focus on supervised learning ds.
unsupervised \rightarrow DPR

1.11.1

What is deep learning?

1.12 \rightarrow Conventional ML is in DPR

1.12.1 \rightarrow What is a neural network? (NN)



14.04.2022

ch 2 Tools for Deep Learning

2.1 Software

Scipy, numpy, Tensorflow

2.2 Hardware

Use Google colab for assignments

2.3 Datasets

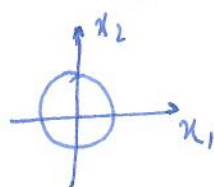
Cityscapes ds \rightarrow why ~~can't~~ wasn't self-supervised learning not performing well?

3. ML basics

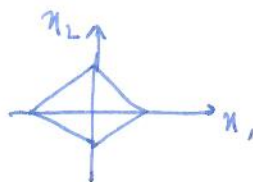
3.1 Linear algebra

See AM for more details

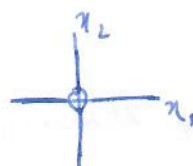
Mainly we use 1. norm in ML



$$p=2: x_1^2 + x_2^2 = 1$$

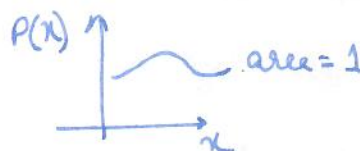
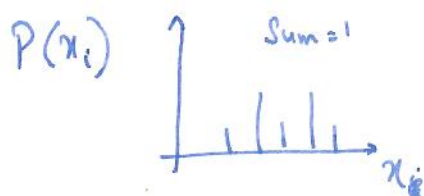


$$p=1: |x_1| + |x_2| = 1$$



$$p=0: \text{one of } x_1, x_2 \neq 0$$

3.2 Random variables and probability distributions.



3.7.1
PDF for discrete valued RVs:

$$p(\underline{x}) = \sum_i p_i \delta(\underline{x} - \underline{x}_i),$$

$\delta(\underline{x})$: dirac f.n.

Cumulative distribution f.n. (CDF)

$$F(\underline{x}) = P(\underline{X} \leq \underline{x}) = \int_{-\infty}^{\underline{x}} p(\underline{z}) d\underline{z}$$

$$p(\underline{x}) = \frac{\partial^d F(\underline{x})}{d x_1 \dots d x_d}$$

3.9.1

Special case $d=1$

$$\underline{X} \rightarrow X \in \mathbb{R}$$

$$\underline{\mu} \rightarrow \mu \in \mathbb{R}$$

$$\begin{aligned} \underline{\sigma} &\rightarrow \text{variance of } \underline{X} = \text{Var}(X) = \sigma^2 \\ &= E[(X - \mu)^2] = E[X^2] - \mu^2 \end{aligned}$$

$\sigma = \sqrt{\text{Var}(X)} \geq 0$: standard ~~varia~~ deviation of X

for any f.n. $g(\underline{X})$ of \underline{X} :

$$E[g(X)] = \int g(\underline{x}) p(\underline{x}) d\underline{x} \stackrel{\text{d.v.}}{=} \sum_i g(\underline{x}_i) P(\underline{x}_i)$$

3.12

Multi noulli distⁿ is an extended case of Bernoulli distⁿ

21.4.22

3.14

One-hot coding

1-hot, 0-cold

One-hot coding - used only in classification

3.15.1

Reformulation of categorical distribution by one-hot coding:

$\underline{x} = [x_i] \in \{\underline{e}_1, \dots, \underline{e}_c\}$, all $x_i = 0$ except for one value at i^{th} position as 1.

$$\Rightarrow \text{PMF } P(\underline{X} = \underline{x}) = P(\underline{X})$$

$$= \begin{cases} p_1 & \text{if } \underline{x} = \underline{e}_1 \quad \text{or } x_1 = 1 \\ \vdots & \vdots \\ p_c & \text{if } \underline{x} = \underline{e}_c \quad \text{or } x_c = 1 \end{cases}$$

$$= p_1^{x_1} \dots p_c^{x_c}$$

$$= \prod_{i=1}^c p_i^{x_i}$$

$$\begin{aligned} \ln P(\underline{X}) &= \sum_{i=1}^c x_i \log p_i = [\underline{x}_1, \dots, \underline{x}_c] \cdot \begin{bmatrix} \ln p_1 \\ \vdots \\ \ln p_c \end{bmatrix} \\ &= \underline{x}^T \underset{\substack{\uparrow \\ \text{element wise}}}{\ln \underline{P}} \end{aligned}$$

3.2.2 Multiple RV

(3)

Product rule $p(\underline{x}, \underline{y}) = p(\underline{x} | \underline{y}) p(\underline{y}) = p(\underline{y} | \underline{x}) p(\underline{x})$

Bayes rule $p(\underline{y} | \underline{x}) = p(\underline{x} | \underline{y}) \cdot \frac{p(\underline{y})}{p(\underline{x})}$

3.17 Chain rule of Probability

$$p(\underline{x}_1, \dots, \underline{x}_N) = p(\underline{x}_1 | \underbrace{\underline{x}_2, \dots, \underline{x}_N}_{\underline{y}}) p(\underline{x}_2, \dots, \underline{x}_N)$$

$\underline{y} \leftarrow$ it can be assumed

Two RV \underline{X} and \underline{Y} independent if

$$p(\underline{x}, \underline{y}) = p(\underline{x}) p(\underline{y})$$

$$\Rightarrow p(\underline{x} | \underline{y}) = p(\underline{x}), \quad p(\underline{y} | \underline{x}) = p(\underline{y})$$

$\underline{x}_1, \dots, \underline{x}_N$ are independent and identically distributed (iid)

$$\ast p(\underline{x}_1, \dots, \underline{x}_N) = \prod_{i=1}^N p_i(\underline{x}_i) \quad \underline{x}_i \sim p_i(\underline{x}_i)$$

$$\ast p_i(\underline{x}_i) = p(\underline{x}_i) \quad \forall i$$

$$\Rightarrow p(\underline{x}_1, \dots, \underline{x}_N) = \prod_{i=1}^N p(\underline{x}_i)$$

3.2.3 Kernel-based density estimation

PDF $p(\underline{x})$ of $\underline{x} \in \mathbb{R}^d$ unknown. Only iid samples

$\underline{x}(n), 1 \leq n \leq N$ available

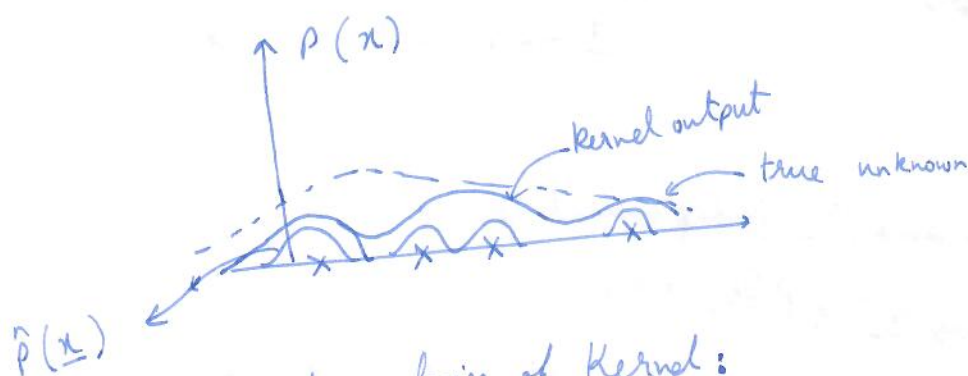
Kernel-based estimate of $p(\underline{x})$ from $\underline{x}(n)$:

Kernel $k(\underline{x})$, like a PDF

$$*) k(\underline{x}) \geq 0 \quad \forall \underline{x}$$

$$*) \int k(\underline{x}) d\underline{x} = 1$$

$$\hat{p}(\underline{x}) = \frac{1}{N} \sum_{n=1}^N k(\underline{x} - \underline{x}(n))$$



popular choice of Kernel:

Gaussian kernel - smooth

$$*) \text{Standard } N(0, \underline{I}) : k(\underline{x}) = \frac{1}{(2\pi)^{d/2}} e^{-\frac{1}{2} \|\underline{x}\|^2}$$

fixed width

$$*) N(0, \sigma^2 \underline{I}) : k(\underline{x}) = \frac{1}{(2\pi\sigma^2)^{d/2}} e^{-\frac{1}{2\sigma^2} \|\underline{x}\|^2}$$

variable ~~time~~ width

Dirac kernel

$$k(\underline{x}) = \delta(\underline{x})$$

$$*) \delta(\underline{x}) = \begin{cases} 0 & \underline{x} \neq 0 \\ \infty & \underline{x} = 0 \end{cases}$$

$$*) \int \delta(\underline{x}) d\underline{x} = 1$$

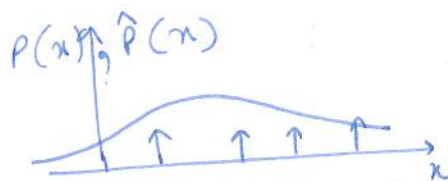
*) Sampling property

$$\int \delta(\underline{x} - \underline{x}_0) f(\underline{x}) d\underline{x} = f(\underline{x}_0)$$

Empirical distribution

(4)

$$\hat{p}(\underline{x}) = \frac{1}{N} \sum_{n=1}^N \delta(\underline{x} - \underline{x}(n))$$



3.13.1

3.3

Kullback - Liebler divergence and cross-entropy

Dissimilarity measure b/w two distribution - KL divergence

A. Continuous-valued RV: PDF

$\underline{X} \sim p(\underline{x})$: true distribution of \underline{X}

$g(\underline{x})$: approximation for $p(\underline{x})$ by a DNN

KL D between p and g :

$$D_{KL}(p \parallel g) = \int p(\underline{x}) \ln \frac{p(\underline{x})}{g(\underline{x})} d\underline{x} \quad \begin{matrix} \geq 0 \\ \leq 0 \end{matrix}$$

So D_{KL} can be +ve, -ve, 0

$$= E \left[\ln \frac{p(\underline{x})}{g(\underline{x})} \right]_{\substack{\underline{x} \sim p \\ \text{expectation over } p(\underline{x})}}$$