

SVEUČILIŠTE U ZAGREBU
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

DIPLOMSKI RAD br. 1728

**Nadzirani pristupi za procjenu
nesigurnosti predikcija dubokih
modela**

Ivan Grubišić

Zagreb, svibanj 2018.

Umjesto ove stranice umetnite izvornik Vašeg rada.

Da bi ste uklonili ovu stranicu obrišite naredbu \izvornik.

Nadzirani pristupi za procjenu nesigurnosti predikcija dubokih modela

Procjena nesigurnosti predikcija vrlo je važan sastojak mnogih praktičnih primjena konvolucijskih modela računalnog vida. Do tog cilja možemo doći analizom višeznačnosti podataka, nesigurnosti odluke modela te vjerojatnosti da se podatak nalazi u distribuciji skupa za učenje. U ovom radu razmatramo pristupe koji procjenu nesigurnosti predikcija uče nadzirano, primjenom istih podataka na kojima se uči i promatrani model.

U okviru rada, potrebno je proučiti i ukratko opisati postojeće pristupe za procjenu nesigurnosti predikcija. Uhodati postupke procjene nesigurnosti dubokih konvolucijskih modela temeljene na nadziranom učenju. Validirati hiperparametre te prikazati i ocijeniti ostvarene rezultate na problemu semantičke segmentacije. Predložiti pravce budućeg razvoja. Radu priložiti izvorni i izvršni kod razvijenih postupaka, ispitne slijedove i rezultate, uz potrebna objašnjenja i dokumentaciju. Citirati korištenu literaturu i navesti dobivenu pomoć.

zahvala

SADRŽAJ

Oznake	vii
1. Uvod	1
2. Osnovni pojmovi	2
2.1. Teorija vjerojatnosti i teorija informacije	2
2.1.1. Slučajne varijable i razdiobe	2
2.1.2. Zdužena, uvjetna i marginalna vjerojatnost i osnovna pravila vjerojatnosti	4
2.1.3. Nezavisnost, uvjetna nezavisnost i uvjetna zavisnost	5
2.1.4. Očekivanje, varijanca i kovarijanca	6
2.1.5. Funkcije slučajnih varijabli	7
2.1.6. Primjeri razdioba	8
2.1.7. Teorija informacije	10
2.2. Nadzirano strojno učenje	13
2.2.1. Induktivna pristranost i komponente algoritma strojnog učenja	13
2.2.2. Kapacitet modela, prenaučenost i podnaučenost	15
2.2.3. Odabir modela	15
2.2.4. Evaluacijske mjere	15
2.3. Procjena parametara i zaključivanje kod probablističkih modela . . .	15
2.3.1. Procjenitelji i točkaste procjene parametara	15
2.3.2. Svojstva i pogreška procjenitelja	16

2.3.3.	Procjenitelj maksimalne izglednosti	16
2.3.4.	Procjenitelj maksimalne aposteriorne vjerojatnosti	17
2.3.5.	Bayesovski procjenitelj i zaključivanje	18
2.4.	Minimizacija rizika	19
2.4.1.	Rizik i empirijski rizik	19
2.4.2.	Strukturni rizik	20
2.5.	Probabilistički grafički modeli	20
2.5.1.	Bayesovski modeli	20
2.6.	Monte Carlo aproksimacija	20
2.7.	Aproksimacija razdioba i aproksimacijsko zaključivanje	21
2.7.1.	Varijacijsko zaključivanje	22
2.7.2.	Monte Carlo aproksimacija	23
3.	Duboko učenje i konvolucijske mreže	24
3.1.	Duboke neuronske mreže	24
3.2.	Konvolucijske mreže	24
3.3.	Optimizacija	24
3.3.1.	Propagacija pogreške unatrag	24
3.3.2.	Isključivanje neurona - dropout	24
3.3.3.	Normalizacija po grupama	24
4.	Procjenjivanje nesigurnosti	25
4.1.	Aleatorna i epistemička nesigurnost	25
5.	Bayesovske neuronske mreže	27
6.	Procjenjivanje nesigurnosti kod konvolucijskih mreža	28
7.	Eksperimentalni rezultati	29
7.1.	Skupovi podataka	29

8. Zaključak	30
Literatura	31
Appendices	33
A. Izvod donje varijacijske granice	33

Oznake

Objekti

Varijable se označavaju kosim slovima sa serifima, većina konstanti uspravnim slovima sa serifima, a slučajne varijable kosim slovima bez serifa. Vektori se označavaju malim podebljanim slovima, matrice i višedimenzionalni nizovi (tenzori) velikim podebljanim slovima, a skupovi slovima s udvostručenim linijama. Za svaku vrstu objekta mogu se koristiti i latinska i grčka slova.

a, A, θ	Varijabla (najčešće skalar)
$\mathbf{a}, \boldsymbol{\theta}$	Vektor ili niz (najčešće vektor stupac)
$\mathbf{A}, \boldsymbol{\Theta}$	Matrica ili višedimenzionalni niz
\mathcal{A}	Skup ili multiskup
$a, A, `$	Konstanta
$\mathbf{a}, `$	Konstanta vektor ili niz
$\mathbf{A}, \boldsymbol{\Theta}$	Konstanta matrica ili višedimenzionalni niz
$\mathbb{A}, \not\mathbb{A}$	Konstanta skup
a, A, θ	Slučajna varijabla
$\mathbf{a}, \boldsymbol{\theta}$	Slučajni vektor ili niz
$\mathbf{A}, \boldsymbol{\Theta}$	Slučajna matrica ili višedimenzionalni niz
\mathcal{A}	Slučajni skup ili multiskup

Konstante

$\{\}$	Prazni skup
$\mathbf{0}$	Nul-vektor
\mathbf{e}_i	i -ti vektor kanonske baze
$\mathbf{1}$	Zbroj svih vektora kanonske baze
\mathbf{I}, \mathbf{I}_n	Matrica identiteta (s n redaka i stupaca)
$\mathbb{N}, \mathbb{Z}, \mathbb{R}, \mathbb{C}$	Poznati skup
$\mathbb{R}_{\geq 0}, \mathbb{R}_{> 0}$	Skup nenegativnih/pozitivnih realnih brojeva

Skupovi i nizovi

$a..b$	Kraći zapis za a, \dots, b
$\{a..b\}$	Skup cijelih brojeva od a do b
$\{a_i: i = 1..n\}, \{a_1..a_n\}, \{a_i\}_{i=1..n}$	Skup s n elemenata

$\{f(a): P(a)\}, \{f(a)\}_{P(a)}, \{f(a)\}_a$	Skup čiji su elementi definirani preko funkcije f i predikata P koji može biti implicitan ili neodređen
$(a_i)_i, (a_{i,j})_{i,j}, (a_{i,j,k})_{i,j,k}$	Višedimenzionalni niz s implicitnim ili neodređenim brojem elemenata
(a, b)	Otvoreni interval
$[a, b]$	Zatvoreni interval

Indeksiranje

Indeksi elemenata vektora ili višedimenzionalnih nizova se radi jednoznačnosti mogu pisati u indeksu oznake vektora u uglatim zagradama. Npr. ako je definiran vektor $\mathbf{a} = (a_1, \dots, a_n)^T$, onda je njegov i -ti element $\mathbf{a}_{[i]} = a_i$.

$\mathbf{a}_{[i]}$	i -ti element vektora \mathbf{a}
$\mathbf{a}_{[i_1:i_2]}$	Vektor kojeg čine elementi $\mathbf{a}_{[i_1]}, \mathbf{a}_{[i_1+1]}, \dots, \mathbf{a}_{[i_2]}$
$\mathbf{a}_{[(i_1 \dots i_n)]}$	Vektor kojeg čine elementi $\mathbf{a}_{[i_1]}, \mathbf{a}_{[i_2]}, \dots, \mathbf{a}_{[i_n]}$
$\mathbf{A}_{[i,j]}$	Element i, j matrice \mathbf{A}
$\mathbf{A}_{[i,:]}$	i -ti redak matrice \mathbf{A}
$\mathbf{A}_{[:,i_1:i_2,j]}$	2-D odsječak 3-D niza \mathbf{A}

Operacije linearne algebre

$\langle \mathbf{a} \mathbf{b} \rangle$	Skalarni produkt, može biti i $\mathbf{a}^T \mathbf{b}$
$\mathbf{a} \odot \mathbf{b}$	Umnožak po elementima; Hadamardov produkt
$\mathbf{a} \oslash \mathbf{b}$	Dijeljenje po elementima
\mathbf{AB}	Matrično množenje
\mathbf{A}^{-1}	Inverz matrice
\mathbf{A}^T	Transponiranje
$\text{diag}(\mathbf{a})$	Dijagonalna matrica kojoj dijagonalu čini vektor \mathbf{a}
$\det \mathbf{A}$	Determinanta matrice \mathbf{A}
$\ \mathbf{a}\ _p$	L^p -norma vektora \mathbf{a}
$\ \mathbf{A}\ _p$	Matrična L^p -norma matrice \mathbf{A}
$\ \mathbf{A}\ _F$	Frobeniusova norma matrice \mathbf{A}

Diferencijalni račun

$\frac{dy}{dx}$	Derivacija y po x
$\frac{\partial y}{\partial x}$	Parcijalna derivacija y po x

$\nabla_x y$	Gradijent y po x
$\nabla_X y$	Gradijent y po X
$\frac{\partial y}{\partial x}, J_{x \mapsto y}, J$	Jakobijan iz $\mathbb{R}^{m \times n}$ za $y \in \mathbb{R}^m$ i $x \in \mathbb{R}^n$
$\int_A f(x) dx, \int_{x \in A} f(x)$	Određeni integral funkcije $f(x)$ po $x \in A$
$\int f(x) dx, \int_x f(x)$	Određeni integral funkcije $f(x)$ po $x \in A$, gdje je A poznat iz konteksta

Teorija vjerojatnosti i teorija informacije

Svakoj slučajnoj varijabli a jednoznačno je dodijeljena jedna razdioba $p(a)$ (ili $P(a)$) i funkcija gustoće vjerojatnosti (poopćena funkcija) $p_a(a) = p(a = a)$. Funkcija gustoće vjerojatnosti se može napisati još na 2 načina. Najkraći zapis je $p(a)$, gdje se po slovu implicitno pretpostavlja slučajna varijabla označena istim slovom bez serifa. $P(A)$ označava vjerojatnost događaja A , a $P_a(a)$ funkciju vjerojatnosti slučajne varijable a . Mogu se koristiti i druge oznake za funkciju vjerojatnosti ili funkciju gustoće vjerojatnosti.

a	Slučajna varijabla
$(a \mid b = b), (a \mid b)$	Uvjetna slučajna varijabla
(a, b)	Združena slučajna varijabla
\mathcal{A}, p, q	Razdioba ili funkcija gustoće vjerojatnosti
A	Događaj
$\{R(a)\}$	Događaj definiran predikatom slučajne varijable a
$P(\{R(a)\}), P(R(a))$	Vjerojatnost događaja $\{R(a)\}$
$P(a), p(a)$	Razdioba slučajne varijable a ; P ako je a diskretna slučajna varijabla, a p ako nije ili ako se ne zna
$P_a(a), P(a)$	Funkcija vjerojatnosti slučajne varijable a : $P_a(a) = P(a = a)$
$p(a = a)$	Gustoća vjerojatnosti događaja $a = a$
$p_a(a), p(a)$	Funkcija gustoće vjerojatnosti slučajne varijable a
$p_{a b}(a), p(a \mid b)$	Gustoća vjerojatnosti za događaj $\{a = a \mid b = b\}$; $p_{a b}(a) = p(a = a \mid b = b)$
$p_{a,b}(a, b), p(a, b)$	Gustoća vjerojatnosti za događaj $\{a = a, b = b\}$; $p_{a,b}(a, b) = p(a = a, b = b)$
$a \sim \mathcal{A}, p(a) = \mathcal{A}$	Slučajna varijabla a ima razdiobu \mathcal{A}
$a \sim A$	Slučajna varijabla a ima takvu razdiobu da svi elementi (multi)skupa A imaju vjerojatnost proporcionalnu višestrukosti $(\frac{1}{ A })$ za običan skup)

$a \sim \mathcal{A}$	a se izvlači iz razidobe \mathcal{A}
$a \sim a, a \sim p(a)$	a se izvlači iz razidobe $p(a)$
$\mathbf{E}_{a \sim a} f(a), \mathbf{E}_a f(a)$	Očekivanje funkcije slučajne varijable a
$\mathbf{D}_{a \sim a} f(a), \mathbf{D}_a f(a)$	Disperzija (varijanca) funkcije slučajne varijable a
$\text{Cov}(a, b)$	Kovarianca
$\mathcal{N}(\mu, \sigma^2)$	Normalna razdioba s učekivanjem μ i varijancom σ^2
$\mathbf{I}(\mathcal{A})$	Sadržaj informacije događaja \mathcal{A}
$H(a)$	Shannonova entropija
$H_b(a)$	Unakrsna entropija
$h(a)$	Diferencijalna entropija
$D_{\text{KL}}(a \parallel b)$	Kullback-Leiblerova divergencija

Ostale oznake

$f: \mathcal{A} \rightarrow \mathcal{B}$	Funkcija s domenom \mathcal{A} i kodomenom \mathcal{B}
$x \mapsto g(x)$	Definicija funkcije; funkcija koja preslikava x iz domene u $g(x)$ iz kodomene
$f * g$	Konvolucija funkcija f i g
$ \mathcal{A} $	Kardinalitet skupa \mathcal{A}
$\delta(\cdot)$	Diracova delta
$\llbracket \cdot \rrbracket$	Iversonova uglata zagrada; $\llbracket P \rrbracket = \begin{cases} 1, & P \equiv \top \\ 0, & P \equiv \perp \end{cases}$

1. Uvod

Uvod rada. Nakon uvoda dolaze poglavlja u kojima se obrađuje tema.

duboko učenje

neizvjesnost modela

primjene procjene nesigurnosti

primjena na semantičkoj segmentaciji i procjeni dubine

struktura rada

2. Osnovni pojmovi

2.1. Teorija vjerojatnosti i teorija informacije

Jako važan pojam u strojnom učenju je nesigurnost ili neizvjesnost. Ona dolazi od šuma u mjerenju i iz konačnosti skupa podataka (Bishop, 2006). Teorija vjerojatnosti nam omogućuje modeliranje nesigurnosti pronalaženje optimalnih zaključaka korištenjem dostupnih informacija.

Postoje dvije glavne interpretacije vjerojatnosti (Murphy, 2012). Jedna je **frekventistička interpretacija** prema kojoj vjerojatnosti predstavljaju učestalosti različitih događaja ako se pokus ponavlja velik broj puta. Druga je **bayesovska interpretacija** prema kojoj vjerojatnost izražava našu nesigurnost o ishodu pokusa.

Ovaj odjeljak daje kratak i matematički ne potpuno precizan pregled nekih od osnovnih pojmova i pravila vezanih uz vjerojatnost. Na strukturu ovog odjeljka imaju utjecaj Goodfellow et al. (2016); Murphy (2012).

2.1.1. Slučajne varijable i razdiobe

Neizvjesnost neke pojave modeliramo **slučajnom varijablom**. Slučajnoj varijabli dodijeljena je **razdioba** koja definira skup vrijednosti koje slučajna varijabla može poprimiti i vjerojatnosti ostvarivanja tih vrijednosti. Skup mogućih vrijednosti neke slučajne varijable još se naziva i **prostor elementarnih događaja**. **Elementarni događaj** je element prostora elementarnih događaja i, ako je x slučajna varijabla za koju se u nekom eksperimentu opaža vrijednost x , taj događaj ima zapis $\{x = x\}$, a njegova vjerojatnost $P(\{x = x\})$ ili $P(x = x)$. **Događaj** je skup vrijednosti i obično se izražava predikatom nad slučajnom varijablom: $\{R(x)\} = \{x: R(x)\}$. Ako je \mathbb{X}

prostor elementarnih događaja slučajne varijable x , onda $P(x \in \mathbb{X}) = 1$. Funkcija

$$\begin{aligned} P_x: \mathbb{X} &\rightarrow [0, 1] \\ x &\mapsto P(x = x) \end{aligned}$$

je **funkcija vjerojatnosti** (engl. *probability mass function, pmf*).

Razlikujemo diskretne i kontinuirane slučajne varijable. Prostor elementarnih događaja diskretne slučajne varijable je prebrojiv skup. Razdioba kontinuirane slučajne varijable x koja poprima vrijednosti iz skupa \mathbb{X} je određena **funkcijom gustoće vjerojatnosti** (engl. *probability density function, pdf*)

$$\begin{aligned} p_x: \mathbb{X} &\rightarrow [0, \infty) \\ x &\mapsto p(x) \end{aligned}$$

za koju vrijedi

$$P(x \in A) = \int_A p_x(x) dx \quad (2.1)$$

za svaki $A \subset \mathbb{X}$.

Funkciju gustoće vjerojatnosti možemo smatrati i **poopćenom funkcijom**¹. To nam omogućuje da funkcijom gustoće predstavljamo razdiobe za koje neki elementarni događaji imaju vjerojatnost veću od 0. Diskretnu razdiobu onda možda možemo predstaviti funkcijom gustoće vjerojatnosti

$$p_x(x) = \sum_{x' \in \mathbb{X}} P(x = x') \delta(x - x'), \quad (2.2)$$

gdje je \mathbb{X} prostor elementarnih događaja slučajne varijable x , a δ Diracova delta, poopćena funkcija za koju vrijedi $\delta(x) = 0$ za $x \neq 0$ i $\int_x \delta(x) dx = 1$. Diracova delta se može promatrati kao limes funkcije gustoće Gaussove razdiobe:

$$\delta(x) = \lim_{\sigma \rightarrow 0} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2}\right).$$

¹[https://en.wikipedia.org/wiki/Distribution_\(mathematics\)](https://en.wikipedia.org/wiki/Distribution_(mathematics))

Ako je x vektor $x = (x_1, \dots, x_n)$, mora vrijediti

$$\delta(x) := \prod_i \delta(x_i). \quad (2.3)$$

Onda n -struki integral gustoće definirane izrazom (2.2) ima vrijednost 1.

Razdioba slučajne varijable x će se u ovom radu označavati s $P(x)$ ako je diskretna, a s $p(x)$ ako je kontinuirana ili neodređena. Funkcija (gustoće) vjerojatnosti će se označavati bez oznake slučajne varijable u indeksu ako je po slovu vrijednosti jasno o kojoj se varijabli radi. Druge oznake koje se koriste opisane su u popisu oznaka na početku rada.

2.1.2. Združena, uvjetna i marginalna vjerojatnost i osnovna pravila vjerojatnosti

Dvije razdiobe su iste ako imaju iste funkcije gustoće vjerojatnosti. Dvije slučajne varijable, i ako imaju istu razdiobu, ne moraju biti iste jer se mogu razlikovati po odnosima s drugim slučajnim varijablama.

Možemo razmatrati više slučajnih varijable zajedno (združenu slučajnu varijablu) i njihovu **združenu razdiobu** $p(x, y)$. Događaji onda imaju oblik $\{R(x, y)\}$. Elementarni događaj onda ima oblik $\{x = x, y = y\}$. Dalje će se izrazi pravila vjerojatnosti odnositi samo na elementarne događaje. Npr. x, y će skraćeno označavati $\{x = x, y = y\}$ kada je jasno po slovima o kojim se slučajnim varijablama radi. Ista pravila vjerojatnosti vrijede i za općenitije događaje jer za svaki događaj možemo definirati indikatorsku slučajnu varijablu kojoj je taj događaj elementarni događaj: $e_i = \llbracket R_i(x, y) \rrbracket$. Takve slučajne varijable imaju skup elementarnih događaja $\{0, 1\}$ i za njih vrijede ista pravila.

Uvjetna vjerojatnost je vjerojatnost nekog događaja ako je poznato da se neki drugi događaj ostvario. Ovako je definirana uvjetna vjerojatnost događaja $\{x = x\}$ ako je poznato da se ostvario događaj $\{y = y\}$:

$$p(x | y) := \frac{p(x, y)}{p(y)}. \quad (2.4)$$

Združena vjerojatnost se može rastaviti **pravilom umnoška**:

$$p(x, y) = p(x | y)p(y). \quad (2.5)$$

Općenitije, pravilo umnoška za n slučajnih varijabli x_1, \dots, x_n izgleda ovako:

$$p(x_1, \dots, x_n) = p(x_1)p(x_2 | x_1) \cdots p(x_n | x_1, \dots, x_{n-1}) \quad (2.6)$$

$$= p(x_1) \prod_{i=2..n} p(x_i | x_1, \dots, x_{i-1}). \quad (2.7)$$

Marginalna vjerojatnost slučajne varijable x je $p(x) = p(x = x, y \in \mathbb{Y})$, gdje je \mathbb{Y} prostor elementarnih događaja slučajne varijable y . Izraženo gustoćom vjerojatnosti (**pravilo zbroja**):

$$p(x) = \int_{\mathbb{Y}} p(x, y) dy = \int_{\mathbb{Y}} p(x | y)p(y) dy. \quad (2.8)$$

Dvije slučajne varijable koje imaju istu razdiobu ne moraju biti u istom odnosu prema drugim slučajnim varijablama. Npr. ako $x_1 \sim \mathcal{A}$, $x_2 \sim \mathcal{A}$ i $y \sim \mathcal{B}$, ne mora vrijediti $p(x_1, y) = p(x_2, y)$.

Rastavljanjem lijeve strane jednadžbe (2.6) na umnožak $p(x | y)p(y)$ dobivamo **Bayesovo pravilo**:

$$p(x | y) = \frac{p(y | x)p(x)}{p(y)}, \quad (2.9)$$

što možemo i ovako zapisati:

$$p(x | y) = \frac{p(y | x)p(x)}{\int p(y | x)p(x) dx}, \quad (2.10)$$

gdje se nazivnik integrira po svim vrijednostima.

2.1.3. Nezavisnost, uvjetna nezavisnost i uvjetna zavisnost

Kada su dvije slučajne varijable x i y **zavisne**, znanje o ishodu jedne utječe na znanje o ishodu druge, tj. uvjetna razdioba $p(x | y = y)$ ovisi o ishodu y . Slučajne varijable x i y su **nezavisne**, što se označava $x \perp y$, akko za svaki par mogućih

vrijednosti (x, y) vrijedi

$$p(x, y) = p(x)p(y), \quad (2.11)$$

ili, ekvivalentno, $p(x | y) = p(x)$ i $p(y | x) = p(y)$. Znanje o ishodu jedne slučajne varijable onda ne utječe na znanje o ishodu druge.

Slučajne varijable x i y , koje mogu biti zavisne, su uz ostvarenje slučajne varijable z **uvjetno nezavisne**, što se označava $x \perp y | z$, akko su slučajne varijable $(x | z = z)$ i $(y | z = z)$ nezavisne za svaki mogući ishod z . To se preko (gustoće) vjerojatnosti može ovako izraziti:

$$p(x, y | z) = p(x | z)p(y | z), \quad (2.12)$$

Isto tako, slučajne varijable x i y koje su nezavisne mogu biti **uvjetno zavisne** uz ostvarenje neke slučajne varijable z . Općenito, dvije slučajne varijable ne moraju biti ni uvjetno zavisne ni uvjetno nezavisne jer neki ishodi treće slučajne varijable mogu utjecati na njihovu zavisnost, a neki ne.

2.1.4. Očekivanje, varijanca i kovarijanca

Očekivanje (prvi moment) slučajne varijable definirano je ovako:

$$\mathbf{E} x := \int x p(x) dx, \quad (2.13)$$

gdje se integrira po prostoru elementarnih događaja. Još se označava ovako: μ_x .

Očekivanje funkcije slučajne varijable zapisujemo ovako:

$$\mathbf{E}_{x \sim x} f(x) := \mathbf{E} f(x) = \int f(x) p(x) dx. \quad (2.14)$$

Ako je po oznaci jasno o kojoj se slučajnoj varijabli radi, možemo kraće pisati

$\mathbf{E}_x f(x)$. Očekivanje ima svojstvo linearnosti:

$$\mathbf{E}[\alpha f(x) + \beta g(x)] = \alpha \mathbf{E} f(x) + \beta \mathbf{E} g(x). \quad (2.15)$$

Varijanca (disperzija, drugi centralni moment) slučajne varijable definirana je

ovako:

$$\mathbf{D}x := \mathbf{E}[(x - \mathbf{E}x)^2] = \int (x - \mathbf{E}x)^2 p(x) dx. \quad (2.16)$$

Varijanca se može izraziti preko drugog momenta $\mathbf{E}x^2$ i kvadrata očekivanja $(\mathbf{E}x)^2$:

$$\mathbf{D}x = \mathbf{E}[(x - \mathbf{E}x)^2] = \mathbf{E}[x^2 - 2x\mathbf{E}x + (\mathbf{E}x)^2] \quad (2.17)$$

$$= \mathbf{E}x^2 - 2(\mathbf{E}x)^2 + (\mathbf{E}x)^2 = \mathbf{E}x^2 - (\mathbf{E}x)^2. \quad (2.18)$$

Drugi korijen varijance je standardna devijacija σ_x .

Kovarijanca para slučajnih varijabli definirana je ovako:

$$\text{Cov}(x, y) := \mathbf{E}[(x - \mathbf{E}x)(y - \mathbf{E}y)] = \mathbf{E}xy - (\mathbf{E}x)(\mathbf{E}y). \quad (2.19)$$

Kovarijacijska matrica slučajnog vektora $\mathbf{x} \in \mathbb{R}^n$ je matrica tipa $n \times n$ takva da:

$$\text{Cov}(\mathbf{x})_{[i,j]} = \text{Cov}(x_{[i]}, x_{[j]}). \quad (2.20)$$

Dijagonalni elementi te matrice su $\text{Cov}(\mathbf{x})_{[i,i]} = \mathbf{D}x_{[i]}$.

2.1.5. Funkcije slučajnih varijabli

Neka je odnos između slučajnih varijabli x i y definiran deterministički funkcijom $f: y = f(x)$. Ako su x i y diskretne slučajne varijable, onda je razdioba slučajne varijable y definirana ovako:

$$P_y(y) = \sum_{x: f(x)=y} P_x(x). \quad (2.21)$$

Ako su x i y kontinuirane slučajne varijable s vrijednostima iz \mathbb{R} i f je injektivna, može se pokazati (Elezović, 2007) da vrijedi

$$p_y(y) = p_x(x) \left| \frac{dx}{dy} \right|. \quad (2.22)$$

To se može poopćiti i na vektore. Onda je $p_y(\mathbf{y}) = \left| \det \frac{\partial \mathbf{x}}{\partial \mathbf{y}} \right|$ (Murphy, 2012).

Neka je z zbroj slučajnih varijabli x i y . Onda vrijedi

$$p_z(z) = \int p_{x,y}(x, z-x) dx. \quad (2.23)$$

Ako su x i y nezavisne, onda to postaje konvolucija:

$$p_z(z) = \int p_x(x)p_y(z-x) dx =: (p_x * p_y)(z). \quad (2.24)$$

2.1.6. Primjeri razdioba

Bernoullijeva razdioba je binarna razdioba s prostorom elementarnih događaja koji je obično $\{0, 1\}$. Ona je onda određena parametrom $\mu \in [0, 1]$ i ima ova svojstva:

$$P(x) = \mu \mathbb{I}[x = 1] + (1 - \mu) \mathbb{I}[x = 0] = \mu^x (1 - \mu)^{1-x}, \quad (2.25)$$

$$\mathbf{E} x = \mu, \quad (2.26)$$

$$\mathbf{D} x = \mu(1 - \mu). \quad (2.27)$$

Kategorička razdioba je poopćenje Bernoullijeve razdiobe na konačan prostor elementarnih događaja koji može imati više od 2 vrijednosti. Ako prostor elementarnih događaja ima kardinalitet n , razdioba je određena vektorom $\mathbf{p} \in [0, 1]^{n-1}$ za koji vrijedi $\sum_i p_{[i]} \leq 1$. Prostor elementarnih događaja ne mora biti skup $\{1..n\}$ pa je kategorička razdioba najopćenitija diskretna razdioba nad konačnim skupom elementarnih događaja.

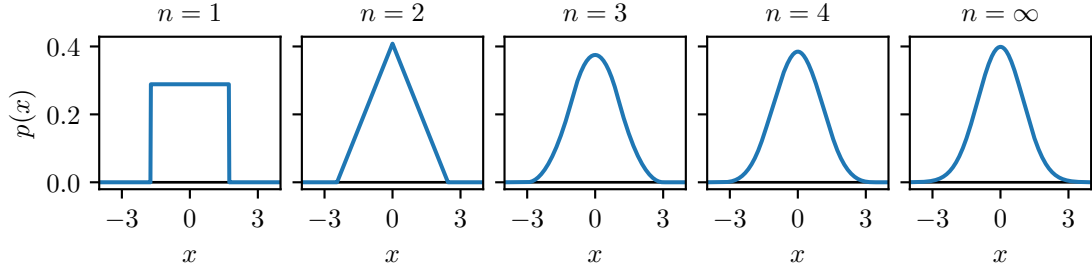
Eksponencijalna razdioba je kontinuirana razdioba s domenom $\mathbb{R}_{\geq 0}$. Ona je definirana parametrom $\lambda \in \mathbb{R}_{>0}$ ili $\beta = \lambda^{-1}$ i ima ova svojstva:

$$p(x) = \lambda \exp(-\lambda x) \quad (2.28)$$

$$\mathbf{E} x = \lambda^{-1}, \quad (2.29)$$

$$\mathbf{D} x = \lambda^{-2}. \quad (2.30)$$

Laplaceova razdioba je kontinuirana razdioba definirana parametrima $\beta \in \mathbb{R}_{>0}$ i



Slika 2.1: Ilustracija centralnog graničnog teorema. Grafovi za različite brojeve pribrojnika n prikazuju funkcije gustoće vjerojatnosti normaliziranih zbrojeva nezavisnih slučajnih varijabli s razdiobom prikazanom prvim grafom. Zadnji graf prikazuje funkciju gustoće Gaussove razdiobe s očekivanjem 0 i varijancom 1.

$\mu \in \mathbb{R}$ i ima ova svojstva:

$$p(x) = \frac{1}{2\beta} \exp\left(-\frac{|x|}{\beta}\right) \quad (2.31)$$

$$\mathbf{E} X = \mu, \quad (2.32)$$

$$\mathbf{D} X = \beta^2. \quad (2.33)$$

Gaussova razdioba ili **normalna razdioba** je kontinuirana razdioba definirana parametrima $\mu \in \mathbb{R}$ i $\sigma \in \mathbb{R}_{>0}$ i ima ova svojstva:

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \quad (2.34)$$

$$\mathbf{E} X = \mu, \quad (2.35)$$

$$\mathbf{D} X = \sigma^2. \quad (2.36)$$

Neka je $z_n = \frac{\sum_{i=1}^n (x_i - \mu)}{\sigma\sqrt{n}}$ normalizirani zbroj n nezavisnih slučajnih varijabli x_i koje imaju jednaku razdiobu s očekivanjem μ i varijancom σ^2 . Prema centralnom graničnom teoremu, z_n u razdiobi konvergira prema Gaussovoj razdiobi kada $n \rightarrow \infty$, tj.

$$\lim_{n \rightarrow \infty} \mathbf{P}(z_n < z) = \int_{-\infty}^z p_{\mathcal{N}(0,1)}(z') dz'. \quad (2.37)$$

To je detaljnije objašnjeno i dokazano npr. u (Elezović, 2007). Centralni granični teorem je ilustriran na slici 2.1.

2.1.7. Teorija informacije

Jedan od osnovnih pojmova u teoriji informacije je **sadržaj informacije** koji događaj preslikava u nenegativan realni broj:

$$I(x \in A) := \log_b \frac{1}{P(x \in A)} = -\log_b P(x \in A). \quad (2.38)$$

Događaji koji imaju manju vjerojatnost sadrže više informacije. Ako je vjerojatnost nekog događaja 1, njegov sadržaj informacije je 0. b je najčešće 2 ili e .

Sadržaj informacije odgovara minimalnom broju simbola (bitova ako $b = 2$) potrebnih za kodiranje elementarnih događaja prefiksnim kodom za koji je očekivanje duljine poruke minimalno (Olah, 2015). Kod prefiksnog koda nijedna kodna riječ nije prefiks neke druge kodne riječi. Takav kod se može prenositi kao niz združenih kodnih riječi bez posebnog simbola za označavanje granica između kodnih riječi. Donja granica očekivanja duljine poruke kod optimalnog koda naziva se **entropija**:

$$H(x) := \mathbf{E}_x I(x = x) = -\mathbf{E}_x \log_b P(x). \quad (2.39)$$

Ona iskazuje neizvjesnost diskretne slučajne varijable. Entropija će biti 0 ako je vjerojatnost nekog elementarnog događaja 1, a najveća će biti kada svi elementarni događaji imaju istu vjerojatnost: $H(x) = \log_b n$, gdje je n broj elementarnih događaja.

Entropija kontinuirane slučajne varijable je beskonačna. Ako se u izrazu (2.39) vjerojatnost zamijeni gustoćom vjerojatnosti, onda on predstavlja **diferencijalnu entropiju**, jedan od analoga² entropije za kontinuirane varijable koji nema neka od svojstava koja ima entropija.

Kao mjera razlike između dviju razdioba često se koristi **Kullback-Leiblerova divergencija** (KL-divergencija) ili **relativna entropija**:

$$D_{KL}(x \parallel y) := \mathbf{E}_x [I(y = x) - I(x = x)] \quad (2.40)$$

$$= \mathbf{E}_x \log_b \frac{P_x(x)}{P_y(x)} \quad (2.41)$$

Ona je uvijek pozitivna i mjeri koliko simbola više se u prosjeku koristi ako se opaža razdioba $P(x)$, a događaji se kodiraju kodom optimalnim za razdiobu $P(y)$, što se

²https://en.wikipedia.org/wiki/Differential_entropy, https://en.wikipedia.org/wiki/Limiting_density_of_discrete_points

bolje vidi ako se ovako izrazi:

$$D_{\text{KL}}(x \parallel y) = H_y(x) - H(x), \quad (2.42)$$

gdje prvi član na desnoj strani jednadžbe označava **unakrsnu entropiju**:

$$H_y(x) := \mathbf{E}_x \mathbf{I}(y = x) = -\mathbf{E}_x \log_b P_y(x). \quad (2.43)$$

Za unakrsnu entropiju se često koristi oznaka $H(x, y)$, ali ista oznaka se koristi i za entropiju združene slučajne varijable (x, y) . Po uzoru na [Olah \(2015\)](#), ovdje koristimo oznaku $H_y(x)$.

KL-divergencija će biti 0 akko x i y imaju iste razdiobe. Ona, kao ni unakrsna entropija, nije simetrična (slika 2.2), tj. općenito $D_{\text{KL}}(x \parallel y) \neq D_{\text{KL}}(y \parallel x)$ i $H_y(x) \neq H_x(y)$. KL-divergencija je izrazom (2.41) definirana i za kontinuirane slučajne varijable ako se funkcije vjerojatnosti zamijene funkcijama gustoće vjerojatnosti. Ona divergira kada postoji x za koji $P_x(x) > 0$ i $P_y(x) = 0$ ili, u slučaju kontinuiranih razdioba, $p_x(x) > 0$ i $p_y(x) = 0$.

Međusobna informacija je mjera zavisnosti između slučajnih varijabli. Definirana je ovako:

$$I(x; y) := \mathbf{E}_{x,y} \log_b \frac{P_{x,y}(x, y)}{P_x(x)P_y(y)}, \quad (2.44)$$

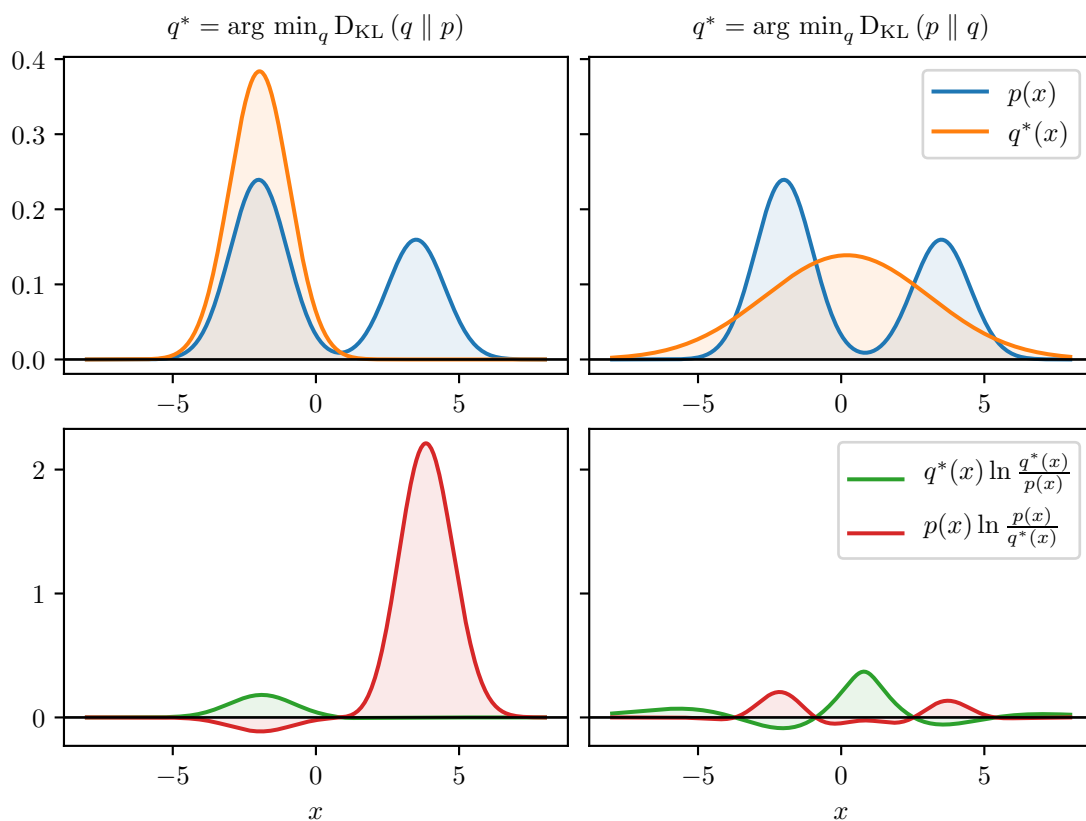
a može se i na ove načine izraziti:

$$I(x; y) = H(x) + H(y) - H(x, y) \quad (2.45)$$

$$= H(x) - H(x | y) \quad (2.46)$$

$$= H(y) - H(y | x). \quad (2.47)$$

Ako su x i y nezavisne, njihova međusobna informacija će biti 0. Ako npr. postoji surjekcija f tako da $y = f(x)$, tj. poznavanje ishoda varijable x jednoznačno određuje ishod varijable y , onda $H(y | x) = 0$ i $I(x; y) = H(y)$. Ako je f bijekcija, onda $I(x; y) = H(x) = H(y)$.



Slika 2.2: Asimetričnost KL-divergencije. p je fiksna razdioba (funkcija gustoće), a q^* je Gaussova razdioba koja ju aproksimira minimizacijom KL-divergencije $D_{\text{KL}}(q \parallel p)$ (lijevo) ili $D_{\text{KL}}(p \parallel q)$ (desno). U donjem retku grafovi prikazuju podintegralne funkcije odgovarajućih KL-divergencija. Kod njih zbrojevi predznačenih površina obojanih područja odgovaraju KL-divergencijama $D_{\text{KL}}(q \parallel p)$ (zeleno) ili $D_{\text{KL}}(p \parallel q)$ (crveno). Optimalna aproksimirajuća razdioba desno ima veliku gustoću gdje god razdioba p ima veliku gustoću. Lijevo optimalna aproksimirajuća razdioba nema veliku gustoću gdje razdioba p nema veliku gustoću. Da je razmak između komponenata razdiobe p malo manji, i lijeva razdioba q^* bi pokrila oba moda i bila sličnija desnoj. Slika je napravljena po uzoru na sliku 3.6 u [Goodfellow et al. \(2016\)](#).

2.2. Nadzirano strojno učenje

Zadatak algoritama nadziranog strojnog učenja je preslikavanje ulaznih primjera $\mathbf{x} \in \mathbb{X}$ u izlaze (oznake) $\mathbf{y} \in \mathbb{Y}$ na temelju konačnog skupa označenih primjera $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_i$. Algoritmima strojnog učenja pretražuje se **model** ili **prostor hipoteza** u cilju pronalaska **hipoteze** koja što bolje **generalizira**, tj. osim primjera iz skupa za učenje, dobro preslikava i neviđene ulazne primjere u izlaze.

Neka je $\mathcal{D} = \{\mathbf{d}_i\}_i$ skup nezavisnih primjera izvučenih iz neke razdiobe \mathcal{D} . Možemo definirati **probabilistički model** \mathcal{H} s nepoznatim parametrima θ kojemu je cilj što bolje modelirati tu razdiobu pronalaskom najbolje hipoteze na temelju podataka: $p(\mathbf{d} \mid \mathcal{D}, \mathcal{H})$. Model koji modelira razdiobu primjera nazivamo **generativnim modelom**. U nastavku ćemo izostavljati oznaku modela radi kraćeg zapisa.

Ako su primjeri parovi $\mathbf{d}_i = (\mathbf{x}_i, \mathbf{y}_i) \in \mathbb{X} \times \mathbb{Y}$, može nam biti cilj ulaznim primjerima iz \mathbb{X} dodjeljivati oznake iz \mathbb{Y} . Ako je problem koji rješavamo dodjeljivanje oznaka ulaznim primjerima, onda su često prikladniji **diskriminativni modeli**. Probabilistički diskriminativni modeli izravno modeliraju uvjetne razdiobe $p(\mathbf{y} \mid \mathbf{x})$ hipotezom koja ulazni primjer \mathbf{x} preslikava u razdiobu $p(\mathbf{y} \mid \mathbf{x}, \mathcal{D})$. Neprobabilistički diskriminativni modeli modeliraju funkciju dodjeljivanja oznaka $f: \mathbb{X} \rightarrow \mathbb{Y}$ hipotezom $h(\mathbf{x})$. Modeliranje zajedničke razdiobe $p(\mathbf{x}, \mathbf{y})$ obično zahtijeva više računalnih resursa i podataka (Bishop, 2006).

2.2.1. Induktivna pristranost i komponente algoritma strojnog učenja

Uz zadani skup hipoteza koji dopušta model, **algoritam strojnog učenja** traži parametre koji definiraju jednu hipotezu. Učenje hipoteze je loše definiran (engl. *ill-posed*) problem jer skup podataka \mathcal{D} nije dovoljan za jednoznačan odabir hipoteze. Osim dobrog opisivanja podataka za učenje, naučena hipoteza mora dobro generalizirati. Kako bi učenje i generalizacija bili mogući, potreban je skup pretpostavki koji se naziva induktivna pristranost. Razlikujemo dvije vrste induktivne pristranosti (Šnajder i Dalbelo Bašić, 2014):

1. **pristranost ograničavanjem** ili **pristranost jezika** – ograničavanje skupa hipoteza koje se mogu prikazati modelom,

2. **pristranost preferencijom** ili **pristranost pretraživanja** – dodjeljivanje različitih prednosti različitim hipotezama.

Većina algoritama strojnog učenja kombinira obje vrste induktivne pristranosti (Šnajder i Dalbelo Bašić, 2014).

Kod većine algoritama strojnog učenja možemo razlikovati 3 osnovne komponente (Šnajder i Dalbelo Bašić, 2014), od kojih prva predstavlja pristranost ograničavanjem, a druge dvije obično pristranost preferencijom:

1. **Model** ili prostor hipoteza. Model \mathcal{H} je skup funkcija h parametriziranih parametrima θ : $\mathcal{H} = \{h(\mathbf{x}; \theta)\}_{\theta}$.
2. **Funkcija pogreške** ili ciljna funkcija. Funkcija pogreške $E(\theta, \mathcal{D})$ na temelju parametara modela (hipoteze) i skupa podataka izračunava broj koji izražava procjenu dobrote hipoteze. Obično pretpostavljamo da su primjeri iz skupa za učenje nezavisni i definiramo **funkcija gubitka** $L: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$, kojoj je prvi parametar predikcija hipoteze, a drugi ciljna oznaka koja odgovara ulaznom primjeru. Funkciju pogreške možemo definirati kao prosječni gubitak na skupu za učenje:

$$E(\theta, \mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} L(h(\mathbf{x}; \theta), \mathbf{y}). \quad (2.48)$$

Obično joj dodajemo **regularizacijski** član kojim unosimo dodatne pretpostavke radi postizanja bolje generalizacije. Više o funkciji pogreške u smislu smanjivanja empirijskog i strukturnog rizika piše u odjeljku 2.4.

3. **Optimizacijski postupak**. Optimizacijski postupak je algoritam kojim pronalazimo hipotezu koja minimizira pogrešku:

$$\theta^* = \arg \min_{\theta} E(\theta, \mathcal{D}). \quad (2.49)$$

Kod nekih jednostavnijih modela minimum možemo odrediti analitički. Inače moramo koristiti neki iterativni optimizacijski postupak. Kod nekih složenijih modela, kao što su neuronske mreže, funkcija pogreške nije unimodalna i vjerojatnost pronalaska globalnog optimuma je zanemariva, ali ipak se mogu pronaći dobra rješenja.

U literaturi riječ *model* često ima šire značenje. Uz skup hipoteza često obuhvaća

i induktivnu pristranost ili dio nje. Riječ *model* može imati i značenje statističkog modela.

2.2.2. Kapacitet modela, prenaučenost i podnaučenost

TODO

2.2.3. Odabir modela

TODO

Murray i Ghahramani (2005)

2.2.4. Evaluacijske mjere

klasifikacija, odnos mF1 i mIoU

2.3. Procjena parametara i zaključivanje kod probablističkih modela

2.3.1. Procjenitelji i točkaste procjene parametara

Ovaj pododjeljak se temelji na Elezović (2007).

Neka je x slučajna varijabla koju promatramo i \mathcal{D} njena razdioba s nama nepoznatim parametrom θ . Taj parametar možemo procijeniti na temelju opaženih vrijednosti x_1, \dots, x_n slučajne varijable x , za što definiramo funkciju g koja daje procjenu parametara

$$\hat{\theta} = f(x_1, \dots, x_N). \quad (2.50)$$

Ako kao parametre takve funkcije uzmemo **uzorak**, tj. skup slučajnih varijabli $\mathcal{D} = (x_1, \dots, x_N)$, gdje pretpostavljamo da su x_1, \dots, x_N međusobno nezavisne i imaju istu razdiobu kao x , dobivamo slučajnu varijablu

$$\hat{\theta} = f(\mathcal{D}). \quad (2.51)$$

Takva slučajna varijabla naziva se **statistika**. Ako je θ nepoznati parametar razdiobe $p(x) = \mathcal{D}$, onda kažemo da je ta statistika $\hat{\theta}$ **procjenitelj** parametra θ , a njen ishod $\hat{\theta}$ **procjena** parametra θ .

2.3.2. Svojstva i pogreška procjenitelja

Pristranost procjenitelja $\hat{\theta}$ je definirana izrazom $\mathbf{E} \hat{\theta} - \theta$, gdje je θ stvarna vrijednost parametra koji se procjenjuje. Ona mjeri koliko procjenitelj griješi neovisno o ishodu uzorka. Kažemo da je procjenitelj parametra θ **nepristran** ako vrijedi

$$\mathbf{E} \hat{\theta} = \theta. \quad (2.52)$$

Varijanca procjenitelja $\hat{\theta}$ je definirana izrazom $\mathbf{D} \hat{\theta}$. Ona mjeri koliko procjenitelj griješi ovisno variranju uzorka. Neka N u oznaci \mathcal{D}_N označava veličinu uzorka. Nepristrani procjenitelj $\hat{\theta}$ je **valjan** ako

$$\lim_{N \rightarrow \infty} \mathbf{D} [\hat{\theta}(\mathcal{D}_N)] = 0. \quad (2.53)$$

Može se pokazati da je očekivanje srednje kvadratne pogreške procjenitelja jednaka zbroju njegove varijance i kvadrata njegove pristranosti (Šnajder i Dalbelo Bašić, 2014), tj.

$$\mathbf{E} [(\hat{\theta} - \theta)^2] = \mathbf{D} \hat{\theta} + (\mathbf{E} \hat{\theta} - \theta)^2. \quad (2.54)$$

2.3.3. Procjenitelj maksimalne izglednosti

Procjenitelj maksimalne izglednosti (ML-procjenitelj, engl. *maximum likelihood*) uzorku dodjeljuje parametre maksimiziraju vjerojatnost uzorka, tj. imaju najveću **izglednost**:

$$\theta_{\text{ML}} = \arg \max_{\theta} p(\mathcal{D} | \theta). \quad (2.55)$$

Zbog pretpostavke međusobne nezavisnosti primjera vrijedi

$$p(\mathcal{D} | \theta) = \prod_{d \in \mathcal{D}} p(d | \theta). \quad (2.56)$$

Za razliku od generativnih, diskriminativni modeli ne modeliraju razdiobu ulaznih primjera, nego samo uvjetnu razdiobu $p(\mathbf{y} \mid \mathbf{x}, \mathcal{D})$ pa kod njih razdioba ulaznih primjera ne ovisi o $\boldsymbol{\theta}$, tj. $p(\mathbf{x} \mid \boldsymbol{\theta}) = p(\mathbf{x})$. Onda je izglednost

$$p(\mathcal{D} \mid \boldsymbol{\theta}) = \prod_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta}) p(\mathbf{x} \mid \boldsymbol{\theta}) = p(\mathbf{x}) \prod_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta}). \quad (2.57)$$

Faktor $p(\mathbf{x})$ ne ovisi o parametrima i može se zanemariti pri optimizaciji.

2.3.4. Procjenitelj maksimalne aposteriorne vjerojatnosti

Procjenitelj maksimalne aposteriorne vjerojatnosti (MAP-procjenitelj, engl. *maximum a posteriori estimator*) u obzir uzima **apriornu razdiobu** $p(\boldsymbol{\theta})$ koja predstavlja dodatne pretpostavke za razdiobu parametara. Apriorna razdioba parametara pojednostavljuje model dajući prednost nekim hipotezama i posebno je korisna kada ima malo podataka. Po Bayesovom pravilu, **aposteriorna vjerojatnost** parametara je

$$p(\boldsymbol{\theta} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{D})} = \frac{p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})}{\int p(\mathcal{D} \mid \boldsymbol{\theta}')p(\boldsymbol{\theta}') d\boldsymbol{\theta}'}. \quad (2.58)$$

Maksimizacijom aposteriorne vjerojatnosti dobivaju se parametri

$$\boldsymbol{\theta}_{\text{MAP}} = \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta} \mid \mathcal{D}) = \arg \max_{\boldsymbol{\theta}} p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta}). \quad (2.59)$$

Ovdje nije potrebno normalizirati aposteriornu vjerojatnost izračunavanjem **marginalne izglednosti** (engl. *marginal likelihood, evidence*) $p(\mathcal{D})$ u nazivniku na desnoj strani jednadžbe (2.58) jer ona ne ovisi o $\boldsymbol{\theta}$, nego samo o modelu \mathcal{H} . Odabirom uniformne apriorne razdiobe MAP-procjenitelj postaje ekvivalentan ML-procjenitelju.

Poželjno je da $p(\mathcal{D} \mid \boldsymbol{\theta})$ i $p(\boldsymbol{\theta})$ kao funkcije parametra $\boldsymbol{\theta}$ imaju takav oblik da njihov umnožak ima sličan oblik i može se analitički izračunati. Ako $p(\boldsymbol{\theta})$ i $p(\boldsymbol{\theta} \mid \mathcal{D})$ imaju isti algebarski oblik definiran nekim parametrima, nazivaju se **konjugatne razdiobe** (Šnajder i Dalbelo Bašić, 2014).

2.3.5. Bayesovski procjenitelj i zaključivanje

Prethodno opisani procjenitelji daju točkastu procjenu parametara i ne izražavaju nesigurnost procjene kojoj uzrok može biti npr. nedovoljna količina podataka ili šum u podacima za učenje. **bayesovski procjenitelj** kao procjenu daje razdiobu nad hipotezama $p(\boldsymbol{\theta} \mid \mathcal{D})$ za koju je integriranjem po svim mogućim parametrima potrebno izračunati marginalnu izglednost $p(\mathcal{D}) = \int p(\mathcal{D} \mid \boldsymbol{\theta}')p(\boldsymbol{\theta}') d\boldsymbol{\theta}'$ iz nazivnika na desnoj strani jednadžbe (2.72).

Kod složenijih modela često ne možemo odabrati konjugatnu apriornu razdiobu, a i funkcija izglednosti je sama po sebi već dovoljno složena da se, neovisno o apriornoj razdiobi, marginalna izglednost $p(\mathcal{D})$ ne može ni analitički ni numerički traktabilno računati.

Vjerojatnost nekog primjera \mathbf{d} procjenjuje se marginalizacijom po svim mogućim parametrima (Neal, 1995):

$$p(\mathbf{d} \mid \mathcal{D}) = \int p(\mathbf{d} \mid \boldsymbol{\theta})p(\boldsymbol{\theta} \mid \mathcal{D}) d\boldsymbol{\theta} = \mathbf{E}_{\boldsymbol{\theta} \mid \mathcal{D}} p(\mathbf{d} \mid \boldsymbol{\theta}). \quad (2.60)$$

Kada se parametri točkasto procjenjuju, npr. MAP-procjeniteljem, točkasta procjena parametara $\hat{\boldsymbol{\theta}}$ aproksimira cijelu aposteriornu razdiobu, tj. $p(\boldsymbol{\theta} \mid \mathcal{D}) \approx \delta(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})$. Onda je

$$p(\mathbf{d} \mid \mathcal{D}) \approx \int p(\mathbf{d} \mid \boldsymbol{\theta})\delta(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) d\boldsymbol{\theta} = p(\mathbf{d} \mid \hat{\boldsymbol{\theta}}). \quad (2.61)$$

Za diskriminativne modele se bayesovsko zaključivanje može izraziti ovako:

$$\begin{aligned} p(\mathbf{y} \mid \mathbf{x}, \mathcal{D}) &= \frac{p(\mathbf{x}, \mathbf{y} \mid \mathcal{D})}{p(\mathbf{x} \mid \mathcal{D})} \\ &= \frac{\int p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta})p(\mathbf{x} \mid \boldsymbol{\theta})p(\boldsymbol{\theta} \mid \mathcal{D}) d\boldsymbol{\theta}}{\int p(\mathbf{x} \mid \boldsymbol{\theta})p(\boldsymbol{\theta} \mid \mathcal{D}) d\boldsymbol{\theta}} \\ &= \frac{p(\mathbf{x}) \int p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta})p(\boldsymbol{\theta} \mid \mathcal{D}) d\boldsymbol{\theta}}{p(\mathbf{x}) \int p(\boldsymbol{\theta} \mid \mathcal{D}) d\boldsymbol{\theta}}. \end{aligned}$$

Poništavanjem $p(\mathbf{x})$ i integriranjem nazivnika dobiva se

$$p(\mathbf{y} \mid \mathbf{x}, \mathcal{D}) = \int p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta})p(\boldsymbol{\theta} \mid \mathcal{D}) d\boldsymbol{\theta} = \mathbf{E}_{\boldsymbol{\theta} \mid \mathcal{D}} p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta}). \quad (2.62)$$

Kod regresije je često, ako pretpostavljamo da pogreška izlaza ima Gaussovu

razdiobu, najbolja procjena hipoteze očekivanje po naučenoj razdiobi parametara (Neal, 1995):

$$h(\mathbf{x}) = \mathbf{E}_{\theta|\mathcal{D}} h(\mathbf{x}; \theta) = \int h(\mathbf{x}; \theta) p(\theta | \mathcal{D}) d\theta. \quad (2.63)$$

U tom slučaju se nesigurnost može izraziti disperzijom $\mathbf{D}_{\theta|\mathcal{D}} h(\mathbf{x}; \theta)$.

2.4. Minimizacija rizika

2.4.1. Rizik i empirijski rizik

Zadatak nadziranog strojnog učenja može se formulirati kao optimizacijski problem minimizacije **rizika**. Neka su θ odabrani parametri. Definiramo **funkciju gubitka** $L: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ koja kažnjava neslaganje izlaza sa stvarnom oznakom. Očekivanje funkcije gubitka je (frekventistički) **rizik** (Murphy, 2012):

$$R(\theta, \mathcal{D}) = \mathbf{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} L(h(\mathbf{x}; \theta), \mathbf{y}). \quad (2.64)$$

Razdioba koja generira podatke nije poznata pa se koristi **empirijski rizik** koji prirodnu razdiobu \mathcal{D} procjenjuje empirijskom, tj. uzorkom \mathcal{D} :

$$R_E(\theta; \mathcal{D}) = \mathbf{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} L(h(\mathbf{x}; \theta), \mathbf{y}) = \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} L(h(\mathbf{x}; \theta), \mathbf{y}). \quad (2.65)$$

U slučaju nenadziranog učenja, kada se hipoteza sastoji od kodera E i dekodera D , tj. $h(\mathbf{x}; \theta) = E(D(\mathbf{x}; \theta); \theta)$, ili generativnog modela, kada je $h(\mathbf{x}; \theta) = p(\mathbf{x} | \theta)$, gubitak mjeri **pogrešku rekonstrukcije** i izraz za rizik je (Murphy, 2012):

$$R(\theta; \mathcal{D}) = \mathbf{E}_{\mathbf{d} \sim \mathcal{D}} L(h(\mathbf{d}; \theta), \mathbf{d}). \quad (2.66)$$

Kod probabilističkih modela empirijski rizik se može definirati kao negativni logaritam izglednosti parametara:

$$R_E(\theta; \mathcal{D}) = -\ln p(\mathcal{D} | \theta) = -\sum_{\mathbf{d} \in \mathcal{D}} \ln p(\mathbf{d} | \theta), \quad (2.67)$$

tj. gubitak je onda $L(h(\mathbf{d}; \theta), \mathbf{d}) = -\ln p(\mathbf{d} | \theta)$. U slučaju diskriminativnog modela, uz zanemarivanja faktora izglednosti koji ne ovisi o θ (jednadžba (2.57)),

vrijedi $L(h(\mathbf{x}; \boldsymbol{\theta}), \mathbf{y}) = -\ln p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta})$.

2.4.2. Strukturni rizik

Kada ima malo podataka ili je model previše složen, minimizacija empirijskog rizika dovodi do velike varijance i slabe generalizacije. Procjenitelj koji minimizira empirijski rizik ne uzima u obzir apriornu razdiobu parametara. Radi postizanja bolje generalizacije, funkciji pogreške dodaje se **regularizacijski** gubitak $\lambda R_R(\boldsymbol{\theta})$, $\lambda \geq 0$, koji predstavlja **strukturni rizik** koji daje prednost jednostavnijim hipotezama:

$$E(\boldsymbol{\theta}; \mathcal{D}) = R_E(\boldsymbol{\theta}; \mathcal{D}) + \lambda R_R(\boldsymbol{\theta}). \quad (2.68)$$

Kod opisanih modela koji uključuju apriorno znanje, regularizacijskom članu uz $\lambda = 1$ odgovara negativni logaritam apriorne vjerojatnosti parametara:

$R_R(\boldsymbol{\theta}) = -\frac{1}{|\mathcal{D}|} \ln p(\boldsymbol{\theta})$. λ različit od 1 odgovara izmjeni entropije apriorne razdiobe $H(p(\boldsymbol{\theta})^\lambda / Z)$, tj. s većim λ apriorna razdioba postaje koncentriranija i regularizacija jača. Jačom regularizacijom se povećava pristranost i smanjuje varijanca procjenitelja.

2.5. Probabilistički grafički modeli

2.5.1. Bayesovski modeli

DL 3.14

2.6. Monte Carlo aproksimacija

Ovaj odjeljak temelji se na [Goodfellow et al. \(2016\)](#).

Monte Carlo aproksimacija je postupak procjenjivanja vrijednosti koje se mogu izraziti kao očekivanje neke funkcije neke slučajne varijable na temelju uzoraka. Ponekad nije moguće analitički ili numerički traktabilno ili efikasno izračunati neki

zbroj ili integral. Ako se on može ovako izraziti:

$$s = \sum_x p(x)f(x) = \mathbf{E} f(x) \quad (2.69)$$

ili

$$s = \int p(x)f(x) dx = \mathbf{E} f(x), \quad (2.70)$$

on se može procijeniti uzorkovanjem:

$$\hat{s}_n = \frac{1}{n} \sum_{i=1}^n f(x_i). \quad (2.71)$$

Procjenitelj \hat{s}_n je nepristran ako su x_i nezavisne i imaju istu razdiobu kao x i valjan ako su varijance $f(x_i)$ ograničene. Vrijedi $\mathbf{D} \hat{s}_n = \frac{1}{n} \mathbf{D} f(x)$.

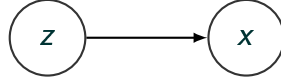
U širem smislu, postupci *Monte Carlo* obuhvaćaju i generiranje uzoraka slučajne varijable čije se očekivanje procjenjuje.

2.7. Aproximacija razdioba i aproksimacijsko zaključivanje

Ovaj odjeljak uglavnom se temelji na [Blei et al. \(2017\)](#) i malo na [Yang \(2017\)](#).

Važan problem u bayesovskoj statistici, gdje se zaključivanje temelji na izračunima koji uključuju aposteriornu razdiobu, je aproksimacija razdioba koje su zahtjevne za računanje. Kod složenijih bayesovskih modela³ aposteriorna razdioba se ne može lako izračunati i treba koristiti aproksimacijske postupke od kojih su glavni **varijacijski** postupci ([Jordan et al., 1999](#)) i **Monte Carlo** postupci koji se temelje na uzorkovanju pomoću **Markovljevog lanca** (MCMC, engl. *Markov chain Monte Carlo*) i Hamiltonovski (ili hibridni) MC-postupci (HMC). MCMC i HMC temelje se na definiranju stohastičkog procesa koji ima stacionarnu razdiobu jednaku razdiobi koja se aproksimira, omogućuju asimptotski egzaktno uzorkovanje i bolje su istraženi. Varijacijski postupci temelje se na aproksimaciji razdiobe nekom jednostavnijom koja se pronalazi rješavanjem optimizacijskog problema, brži su i jednostavniji za ostvariti za složenije modele.

³Bayesovski model (ili Bayesova mreža) je probabilistički grafički model sa strukturom usmjerenog acikličkog grafa.



Slika 2.3: Prikaz grafičkog modela čija je združena razdioba $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{z})p(\mathbf{x} | \mathbf{z})$.

Razmatramo bayesovski model koji ima jednu latentnu varijablu \mathbf{z} i jednu vidljivu varijablu \mathbf{x} . Model je prikazan an slici 2.3 i opisan je ovom jednadžbom združene vjerojatnosti:

$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{z})p(\mathbf{x} | \mathbf{z}).$$

Zaključivanjem se određuje aposteriorna razdioba latentne varijable:

$$p(\mathbf{z} | \mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})} = \frac{p(\mathbf{x}, \mathbf{z})}{\int p(\mathbf{x}, \mathbf{z}) d\mathbf{z}}. \quad (2.72)$$

na temelju opažanih vrijednosti slučajne varijable \mathbf{x} (podataka). Kod složenijih modela integriranje marginalne izglednosti u nazivniku nije traktabilno i aposteriorna razdioba se mora aproksimirati **aproksimacijskim zaključivanjem**.

2.7.1. Varijacijsko zaključivanje

Za razliku od uzorkovanja kod MCMC-postupaka, osnovna ideja kod varijacijskog zaključivanja je optimizacija. Prvo se odabire familija razdioba $\mathcal{Q} = \{p(\tilde{\mathbf{z}})\}_{\tilde{\mathbf{z}}} = \{\mathcal{Q}_\phi\}_\phi$ koje su lakše za računanje. Razdiobe iz \mathcal{Q} su parametrizirane tzv. **varijacijskim parametrima** ϕ . Cilj je na temelju podataka kao zamjenu za aposteriornu razdiobu $p(\mathbf{z} | \mathbf{x})$ pronaći razdiobu iz \mathcal{Q} koja ju što bolje aproksimira. To možemo ostvariti minimizacijom Kullback-Leiblerove (KL) divergenciju s obzirom na stvarnu aposteriornu razdiobu po varijacijskim parametrima:

$$\mathcal{Q}^* = \arg \min_{p(\tilde{\mathbf{z}}) \in \mathcal{Q}} D_{\text{KL}}(\tilde{\mathbf{z}} \parallel (\mathbf{z} | \mathbf{x})). \quad (2.73)$$

Ovakva ciljna funkcija obično se ne može lako izračunati jer zahtijeva računanje

marginalne izglednosti $p(\mathbf{x})$ iz nazivnika u jednadžbi (2.72) marginalizacijom po \mathbf{z} :

$$\begin{aligned} D_{\text{KL}}(\tilde{\mathbf{z}} \parallel (\mathbf{z} \mid \mathbf{x})) &= \mathbb{E}_{\tilde{\mathbf{z}}} \ln \frac{p(\tilde{\mathbf{z}})}{p(\mathbf{z}=\tilde{\mathbf{z}} \mid \mathbf{x})} \\ &= \mathbb{E}_{\tilde{\mathbf{z}}} \ln p(\tilde{\mathbf{z}}) - \mathbb{E}_{\tilde{\mathbf{z}}} \ln p(\mathbf{z}=\tilde{\mathbf{z}} \mid \mathbf{x}) \\ &= \mathbb{E}_{\tilde{\mathbf{z}}} \ln p(\tilde{\mathbf{z}}) - \mathbb{E}_{\tilde{\mathbf{z}}} \ln p(\mathbf{z}=\tilde{\mathbf{z}}, \mathbf{x}) + \ln p(\mathbf{x}). \end{aligned} \quad (2.74)$$

Marginalna izglednost se može zanemariti jer marginalna izglednost ne ovisi o parametrima po kojima se optimizira pa umjesto minimizacije KL-divergencije maksimiziramo funkciju koja se naziva **varijacijska donja granica** (engl. *variational lower bound*) ili **donja granica (logaritma) marginalne izglednosti** (engl. *(log) evidence lower bound, ELBO*):

$$L_x(\tilde{\mathbf{z}}) := \mathbb{E}_{\tilde{\mathbf{z}}} \ln p(\mathbf{z}=\tilde{\mathbf{z}}, \mathbf{x}) - \mathbb{E}_{\tilde{\mathbf{z}}} \ln p(\tilde{\mathbf{z}}). \quad (2.75)$$

Vrijedi $D_{\text{KL}}(\tilde{\mathbf{z}} \parallel (\mathbf{z} \mid \mathbf{x})) = -L_x(\tilde{\mathbf{z}}) + \ln p(\mathbf{x})$. Varijacijska donja granica se može i ovako izraziti:

$$L_x(\tilde{\mathbf{z}}) = \mathbb{E}_{\tilde{\mathbf{z}}} \ln p(\mathbf{x} \mid \mathbf{z}=\tilde{\mathbf{z}}) - D_{\text{KL}}(\tilde{\mathbf{z}} \parallel \mathbf{z}). \quad (2.76)$$

Maksimiziranje takve ciljane funkcije daje razdiobu $p(\tilde{\mathbf{z}})$ koja dobro objašnjava podatke jer se potiče veće očekivanje logaritma izglednosti, i ne razlikuje se previše od apriorne razdiobe jer se potiče manja KL-divergencija između varijacijske razdiobe i apriorne razdiobe (Gal i Ghahramani, 2015).

Naziv *donja granica marginalne izglednosti* dolazi od toga što su Jordan et al. (1999) izveli nejednakost $\ln p(\mathbf{x}) \geq L_x(\tilde{\mathbf{z}})$ preko Jensenove nejednakosti. Ta nejednakost slijedi i iz prethodne jednadžbe i nenegativnosti KL-divergencije:

$$\ln p(\mathbf{x}) = L_x(\tilde{\mathbf{z}}) + D_{\text{KL}}(\tilde{\mathbf{z}} \parallel (\mathbf{z} \mid \mathbf{x})) \geq L_x(\tilde{\mathbf{z}}). \quad (2.77)$$

TODO mean-field, calculus of variations

2.7.2. Monte Carlo aproksimacija

1.3.1 Neal (1995).

3. Duboko učenje i konvolucijske mreže

3.1. Duboke neuronske mreže

3.2. Konvolucijske mreže

3.3. Optimizacija

3.3.1. Propagacija pogreške unatrag

3.3.2. Isključivanje neurona - dropout

3.3.3. Normalizacija po grupama

4. Procjenjivanje nesigurnosti

4.1. Aleatorna i epistemička nesigurnost

Kod bayesovskih modela nesigurnost zaključivanja izražava se razdiobom po vrijednostima varijable čija vrijednost se procjenjuje, a može se izraziti i entropijom ili varijancom kada je prikladno.

Postoje različiti izvori nesigurnosti (C. Kennedy i O'Hagan, 2002), ali nesigurnost općenito možemo podijeliti na dvije vrste (Kiureghian i Ditlevsen, 2009): **aleatornu nesigurnost** i **epistemičku nesigurnost**. Riječ *aleatorna* izvedena je vjerojatno od latinske riječi *aleator* (Gal, 2016) koja znači *kockar*, a riječ *epistemička* izvedena je od grčke riječi *epistēmē* koja znači *znanje*. Aleatorna nesigurnost je nesigurnost koju model ne može smanjiti neovisno o znanju i količini dostupnih podataka. Ona dolazi od nedeterminizma samog procesa koji generira podatke, nedostupnosti dijela informacija ili ograničenja modela. Epistemička nesigurnost je nesigurnost u strukturu i parametre modela Gal (2016). Ona dolazi od neznanja i može se smanjiti uz više podataka.

Granica između aleatorne i epistemičke nesigurnosti ovisi o modelu. Nešto što je kod jednostavnijeg modela aleatorna nesigurnost, kod složenijeg modela može biti će epistemičkog karaktera. Ako su neke pojave po prirodi nasumične ili se ne mogu ili ne žele modelu dati informacije koje bi ih mogle objasniti, nesigurnost zaključivanja u vezi tih pojava će, neovisno o ograničenosti modela, biti aleatorna.

TODO: homoskedastička, heteroskedastička nesigurnost

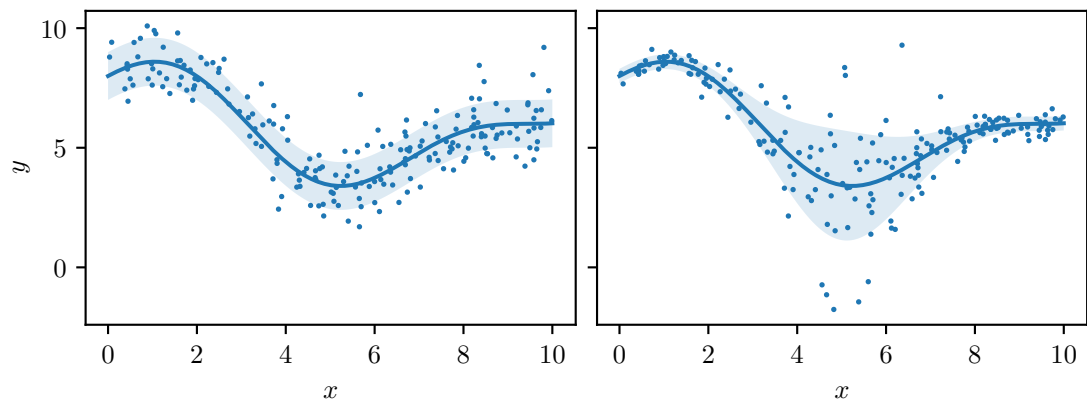
TODO: model uncertainty Gal-thesis 1.2

5.33331pt

11.74983pt

small 10.95pt

footnotesize 10.0pt



Slika 4.1: Homoskedastički (lijevo) i heteroskedastički (desno) Gaussov šum. Crta prikazuje očekivanje $f(x)$, svjetloplava površina standardnu devijaciju šuma $s(x)$, a točke slučajne uzorke. Točke su generirane prema $(y | x) \sim \mathcal{N}(f(x), s(x)^2)$. Na lijevoj slici je $s(x) = 1$.

412.56497pt

5. Bayesovske neuronske mreže

6. Procenjivanje nesigurnosti kod konvolucijskih mreža

7. Eksperimentalni rezultati

7.1. Skupovi podataka

8. Zaključak

Zaključak.

LITERATURA

- Christopher M. Bishop. **Pattern Recognition and Machine Learning (Information Science and Statistics)**. 2006.
- David M. Blei, Alp Kucukelbir, i Jon D. McAuliffe. Variational Inference: A Review for Statisticians. **Journal of the American Statistical Association**, 2017. URL <http://arxiv.org/abs/1601.00670>.
- Marc C. Kennedy i Anthony O'Hagan. Bayesian calibration of computer models. 2002.
- Neven Elezović. **Vjerojatnost i statistika: Slučajne varijable**. 2007.
- Yarin Gal. **Uncertainty in Deep Learning**. Doktorska disertacija, University of Cambridge, 2016.
- Yarin Gal i Zoubin Ghahramani. Dropout as a Bayesian Approximation: Appendix. 2015. URL <https://arxiv.org/abs/1506.02157>.
- Ian Goodfellow, Yoshua Bengio, i Aaron Courville. **Deep Learning**. MIT Press, 2016. <http://www.deeplearningbook.org>.
- Michael I. Jordan, Zoubin Ghahramani, Tommi S. Jaakkola, i Lawrence K. Saul. An introduction to variational methods for graphical models. 1999.
- Armen Der Kiureghian i Ove Ditlevsen. Aleatory or epistemic? Does it matter? 2009.
- Kevin P. Murphy. **Machine Learning: A Probabilistic Perspective**. 2012.
- Iain Murray i Zoubin Ghahramani. A note on the evidence and Bayesian Occam's razor. 2005.
- Jan Šnajder i Bojana Dalbelo Bašić. **Strojno učenje**. 2014.
- Radford M. Neal. Bayesian learning for neural networks, 1995.
- Christopher Olah. Visual information theory, 2015. URL <http://colah.github.io/posts/2015-09-Visual-Information/>.
- Xitong Yang. Understanding the Variational Lower Bound, 2017. URL <http://legacydirls.umiaccs.umd.edu/~xyang35/files/understanding-variational-lower.pdf>.

Nadzirani pristupi za procjenu nesigurnosti predikcija dubokih modela

Sažetak

Sažetak na hrvatskom jeziku.

Ključne riječi: Ključne riječi, odvojene zarezima.

Title

Abstract

Abstract.

Keywords: Keywords.

A. Izvod donje varijacijske granice

The contents...