

SVEUČILIŠTE U ZAGREBU
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

DIPLOMSKI RAD br. 1728

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

Ivan Grubišić

Zagreb, lipanj 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
1.1. Struktura rada	1
2. Osnovni pojmovi	2
2.1. Teorija vjerojatnosti	2
2.1.1. Slučajne varijable i razdiobe	2
2.1.2. Združena, 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
2.2. Teorija informacije	9
2.3. Optimizacija temeljena na gradijentu	13
2.3.1. Gradijentni spust	13
2.3.2. Postupci drugog reda	14
3. Statističko modeliranje	15
3.1. Probabilistički grafički modeli	15
3.2. Procjena parametara i zaključivanje	18
3.2.1. Procjenitelji i točkaste procjene parametara	18

3.2.2.	Svojstva i pogreška procjenitelja	19
3.2.3.	Procjenitelj maksimalne izglednosti	19
3.2.4.	Procjenitelj maksimalne aposteriorne vjerojatnosti	20
3.2.5.	Bayesovski procjenitelj i zaključivanje	20
3.3.	Monte Carlo aproksimacija	22
3.4.	Aproksimacija razdioba i aproksimacijsko zaključivanje	22
3.5.	Postupci uzorkovanja	23
3.6.	Varijacijsko zaključivanje	23
3.6.1.	Metoda polja sredina	25
4.	Nadzirano strojno učenje	26
4.1.	Induktivna pristranost	27
4.2.	Komponente algoritma strojnog učenja	27
4.3.	Kapacitet modela, podnaučenost i prenaučenost	28
4.4.	Odabir modela	29
4.5.	Funkcija pogreške	30
4.5.1.	Rizik i empirijski rizik	30
4.5.2.	Strukturni rizik i regularizacija	31
4.6.	Osnovni zadaci nadziranog učenja	31
4.7.	Primjeri modela: poopćeni linearni modeli	32
5.	Duboko učenje i konvolucijske mreže	35
5.1.	Duboke unaprijedne mreže	36
5.2.	Konvolucijske mreže	38
5.3.	Učenje	38
5.3.1.	Algoritam propagacije pogreške unatrag	38
5.3.2.	Optimizacijski algoritmi	38
5.3.3.	Algoritam propagacija pogreške unatrag	38

5.3.4. Isključivanje neurona - dropout	38
5.3.5. Normalizacija po grupama	38
6. Procenjivanje nesigurnosti	39
6.1. Aleatorna i epistemička nesigurnost	39
6.2. Homoskedastička i heteroskedastička nesigurnost	39
7. Bayesovske neuronske mreže	41
8. Procenjivanje nesigurnosti kod konvolucijskih mreža	42
9. Eksperimentalni rezultati	43
9.1. Programska izvedba	43
9.2. Skupovi podataka	43
10. Zaključak	44
Literatura	45

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}, \mathbb{Z}	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
a , oznaka	Oznaka koja ne predstavlja matematički objekt

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

Definiranje skupova i nizova

$a..b$	Kraći zapis za a, \dots, b
$\{a..b\}$	Skup cijelih brojeva od a do b

$\{f(a): P(a)\}, \{f(a)\}_{P(a)}$	Skup čiji su elementi definirani preko funkcije f i predikata P
$\{f(a)\}_a$	Skup čiji su elementi definirani preko funkcije f i varijabli a iz implicitno određenog skupa
$\{a_1 \dots a_n\}, \{a_i\}_{i=1 \dots n}$	Skup s n elemenata
$(a_i)_i, (a_{i,j})_{i,j}, (a_{i,j,k})_{i,j,k}$	Višedimenzionalni niz s implicitnim ili neodređenim brojem elemenata
(a, b)	Otvoreni interval
$[a, b]$	Zatvoreni interval
$[x_1, \dots, x_n]$	Vektor redak
(x_1, \dots, x_n)	Vektor (stupac)

Donji i gornji indeks

U donjem indeksu oznake mogu biti oznake drugih matematičkih objekata. U donjem i gornjem indeksu oznake mogu biti oznake (slova ili riječi) koje ne predstavljaju matematičke objekte. Oznake koje ne predstavljaju matematičke objekte istog su stila kao tekst. Indeksi (redni brojevi) elemenata vektora ili višedimenzionalnih nizova se, ako nije određeno drugačije, pišu u donjem 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 $\mathbf{a}_{[i]} = a_i$.

a_d^g	Oznaka varijable s oznakama u donjem i gornjem indeksu
$\mathbf{a}_{[i]}$	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 i druge operacije s nizovima

$\langle \mathbf{a} \mathbf{b} \rangle$	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}
$\ \mathbf{a}\ _2$	L^2 -norma vektora \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}
$f(\mathbf{a})$	Ako f nije drugačije definirana i inače označava funkciju $\mathbb{R} \rightarrow \mathbb{R}$, onda se primjenjuje po svakom elementu vektora posebno
$\mathbf{a} \# \mathbf{b}$	Konkatenacija vektora (stupaca) $\mathbf{a} \in \mathbb{R}^n$ i $\mathbf{b} \in \mathbb{R}^m$ u vektor iz \mathbb{R}^{n+m}
$\mathbf{A} \# \mathbf{B}$	Konkatenacija višedimenzionalnih nizova po prvoj dimenziji
$\mathbf{A} \#' \mathbf{B}$	Konkatenacija višedimenzionalnih nizova po zadnjoj dimenziji

Diferencijalni račun

$\frac{dy}{dx}, \frac{d}{dx} f(x)$	Derivacija $y = f(x)$ po x
$\frac{\partial y}{\partial x}, \frac{\partial}{\partial x} f(x)$	Parcijalna derivacija $y = f(x)$ po x
$\nabla_x y, \nabla_x f(x), \left(\frac{\partial y}{\partial x}\right)^T$	Gradijent $y = f(\mathbf{x})$ po \mathbf{x}
$\nabla_X y, \nabla_X f(x)$	Gradijent $y = f(\mathbf{x})$ po \mathbf{X}
$\frac{\partial^2 y}{\partial x \partial x^T}, \mathbf{H}_f(\mathbf{x}), \mathbf{H}$	Hessijan iz $\mathbb{R}^{n \times n}$ za $f: \mathbb{R}^n \rightarrow \mathbb{R}$ i $y = f(\mathbf{x})$
$\frac{\partial \mathbf{y}}{\partial \mathbf{x}}, \mathbf{J}_f(\mathbf{x}), \mathbf{J}$	Jakobijeva matrica iz $\mathbb{R}^{m \times n}$ za $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ i $\mathbf{y} = f(\mathbf{x})$
$\int_A f(x) dx, \int_{x \in A} f(x)$	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 implicitan

Teorija vjerojatnosti

Svakoj slučajnoj varijabli a jednoznačno je dodijeljena jedna razdioba $p(a)$ (ili $P(a)$) i funkcija gustoće vjerojatnosti (koja može biti poopćena funkcija) $p_a(a) = p(a = a)$. $P(A)$ označava vjerojatnost događaja A , a P_a funkciju vjerojatnosti slučajne varijable a . Gustoća vjerojatnosti se još kraće može zapisati $p(a)$, gdje se po slovu implicitno pretpostavlja slučajna varijabla označena istim slovom bez serifa. Isto tako, vjerojatnost elementarnog događaja se može zapisati $P(a)$. Mogu se koristiti i druge oznake za funkciju vjerojatnosti ili funkciju gustoće vjerojatnosti.

$(a \mid b = b), (a \mid b)$	Uvjetna slučajna varijabla
(a, b)	Združena slučajna varijabla
$a \perp b$	<i>Slučajne varijable a i b su nezavisne</i>
$a \not\perp b$	<i>Slučajne varijable a i b su zavisne</i>
$a \perp b \mid c$	<i>Slučajne varijable a i b su uvjetno nezavisne uz poznat ishod slučajne varijable c</i>
$a \not\perp b \mid c$	<i>Slučajne varijable a i b su uvjetno zavisne uz poznat ishod slučajne varijable c</i>
p, q	Razdioba ili funkcija gustoće vjerojatnosti
\mathcal{A}	Događaj
$\{R(a)\}$	Događaj definiran predikatom slučajne varijable a
$P(\{R(a)\}), P(R(a))$	Vjerojatnost događaja $\{R(a)\}$
$P(a), p(a), \mathcal{D}$	Razdioba slučajne varijable a ; P ako je a diskretna slučajna varijabla, p ako nije ili ako se ne zna
$P(a = a), P_a(a), P(a)$	Vjerojatnost događaja $\{a = a\}$
$p(a = a), p_a(a), p(a)$	Gustoća vjerojatnosti događaja $\{a = a\}$
$p_{a b}(a), p(a \mid b)$	Gustoća vjerojatnosti događaja $\{a = a \mid b = b\}$
$p_{a,b}(a, b), p(a, b)$	Gustoća vjerojatnosti događaja $\{a = a, b = b\}$
$a \sim q, p(a) = q$	<i>Slučajna varijabla a ima razdiobu q</i>
$a \sim \mathcal{A}$	<i>Slučajna varijabla a ima takvu razdiobu da svi elementi (multi)skupa \mathcal{A} imaju vjerojatnost proporcionalnu višestrukosti ($\frac{1}{ \mathcal{A} }$ za običan skup)</i>
$a \sim q$	<i>a se izvlači iz razdiobe q</i>
$a \sim a, a \sim p(a)$	<i>a se izvlači iz razdiobe $p(a)$</i>
$\mathbf{E}_{a \sim a} f(a), \mathbf{E}_a f(a)$	Očekivanje funkcije slučajne varijable a
$\mathbf{D}_{a \sim a} f(a), \mathbf{D}_a f(a)$	Disperzija (varijanca) funkcije slučajne varijable a
$\text{Cov}(a, b)$	Kovarijanca
$\mathcal{N}(\mu, \sigma^2)$	Normalna razdioba s očekivanjem μ i varijancom σ^2
$\mathcal{U}(\mathcal{A})$	Uniformna razdioba nad skupom \mathcal{A}

Teorija informacije

$I(\mathcal{A})$	Sadržaj informacije događaja \mathcal{A}
$H(a)$	Entropija

$h(a)$	Diferencijalna entropija
$I(a, b)$	Međusobna informacija
$H(a \mid b)$	Uvjetna entropija
$H_b(a)$	Unakrsna entropija
$D_{\text{KL}}(a \parallel b)$	Kullback-Leiblerova divergencija (relativna entropija)

Grafovi

$\text{pa}_G(a)$	Skup čvorova roditelja čvora a u grafu G
$\text{ch}_G(a)$	Skup čvorova djece čvora a u grafu G
$\text{pred}_G(a)$	Skup čvorova prethodnika čvora a u grafu G
$\text{succ}_G(a)$	Skup čvorova nasljednika čvora a u grafu G

Ostale oznake

$f: A \rightarrow B$	Funkcija s domenom A i kodomenom 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
$ A $	Kardinalitet skupa A
$\delta(\cdot)$	Diracova delta
$\llbracket \cdot \rrbracket$	Iversonova uglatna 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

1.1. Struktura rada

2. Osnovni pojmovi

2.1. Teorija vjerojatnosti

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.

Ovo poglavlje daje kratak i matematički ne potpuno precizan pregled nekih od osnovnih pojmova i pravila vezanih uz vjerojatnost. Na strukturu ovog poglavlja 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

$$P_x: \mathbb{X} \rightarrow [0, 1]$$

$$x \mapsto P(x = x)$$

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*)

$$p_x: \mathbb{X} \rightarrow [0, \infty)$$

$$x \mapsto p(x)$$

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:

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

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

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

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

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

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. Na nekim mjestima će, radi kratkoće, riječ *razdioba* imati značenje *funkcija gustoće* ili *funkcija vjerojatnosti*.

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, marginalizacija**):

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

Dvije slučajne varijable koje imaju istu razdiobu ne moraju biti u istom odnosu prema drugim slučajnim varijablama. Npr. ako $x_1 \sim q_1$, $x_2 \sim q_1$ i $y \sim q_2$, 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**, što se označava $x \not\perp y$, znanje o ishodu jedne utječe na znanje o ishodu druge, tj. uvjetna razdioba $p(x | y = y)$ ovisi o ishodu y . *Znanje o ishodu* ne mora značiti da je ishod poznat. Dovoljna je promjena znanja o razdiobi koja može biti posljedica opažanja neke treće slučajne varijable. Slučajne varijable x i y su **nezavisne**, što se označava $x \perp y$, akko za svaki par (x, y) vrijedi

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

ili, ekvivalentno,

$$p(x | y) = p(x). \quad (2.12)$$

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 znanje o ishodu 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 . Onda za svaku trojku (x, y, z) vrijedi

$$p(x, y | z) = p(x | z) p(y | z), \quad (2.13)$$

ili, ekvivalentno,

$$p(x | y, z) = p(x | z). \quad (2.14)$$

Isto tako, slučajne varijable x i y koje su nezavisne mogu biti **uvjetno zavisne** uz znanje o ishodu 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. Također se može govoriti i o zavisnosti ili nezavisnosti pojedinih događaja.

2.1.4. Očekivanje, varijanca i kovarijanca

Očekivanje (prvi moment) slučajne varijable definirano je ovako:

$$\mathbf{E} x := \int x p(x) dx, \quad (2.15)$$

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

Očekivanje funkcije slučajne varijable zapisujemo ovako:

$$\mathbf{E}_{x \sim x} f(x) := \mathbf{E} f(x) = \int f(x) p(x) dx. \quad (2.16)$$

Ako je po oznaci jasno o kojoj se slučajnoj varijabli radi, možemo kraće pisati $\mathbf{E}_x f(x)$. Očekivanje ima svojstvo linearnosti:

$$\mathbf{E}[\alpha f(x) + \beta g(x)] = \alpha \mathbf{E} f(x) + \beta \mathbf{E} g(x). \quad (2.17)$$

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

ovako:

$$\mathbf{D}x := \mathbf{E}[(x - \mathbf{E}x)^2] = \int (x - \mathbf{E}x)^2 p(x) dx. \quad (2.18)$$

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

$$\mathbf{D}x = \mathbf{E}[(x - \mathbf{E}x)^2] = \mathbf{E}[x^2 - 2x\mathbf{E}x + (\mathbf{E}x)^2] \quad (2.19)$$

$$= \mathbf{E}x^2 - 2(\mathbf{E}x)^2 + (\mathbf{E}x)^2 = \mathbf{E}x^2 - (\mathbf{E}x)^2. \quad (2.20)$$

Drugi korijen varijance je standardna devijacija σ_x .

Kovarijanca para slučajnih varijabli definirana je ovako:

$$\text{Cov}(x, y) := \mathbf{E}[(x - \mathbf{E}x)(y - \mathbf{E}y)] = \mathbf{E}xy - (\mathbf{E}x)(\mathbf{E}y). \quad (2.21)$$

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

$$\text{Cov}(\mathbf{x})_{[i,j]} = \text{Cov}(x_{[i]}, x_{[j]}). \quad (2.22)$$

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 funkcijom f koja ishode jedne slučajne varijable deterministički preslikava u ishode druge, što se označava ovako: $y = f(x)$. Ako su x i y diskretne slučajne varijable, onda je razdioba slučajne varijable y definirana ovako:

$$P_y(y) = \sum_{x: f(x)=y} P_x(x). \quad (2.23)$$

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

$$p_y(y) = p_x(x) \left| \frac{dx}{dy} \right|. \quad (2.24)$$

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

$$p_z(z) = \int p_{x,y}(x, z - x) dx. \quad (2.25)$$

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

$$p_z(z) = \int p_x(x)p_y(z-x)dx =: (p_x * p_y)(z). \quad (2.26)$$

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:

$$P(x) = \mu \mathbb{I}[x = 1] + (1 - \mu) \mathbb{I}[x = 0] = \mu^x (1 - \mu)^{1-x}, \quad (2.27)$$

$$\mathbf{E} x = \mu, \quad (2.28)$$

$$\mathbf{D} x = \mu(1 - \mu). \quad (2.29)$$

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:

$$p(x) = \lambda \exp(-\lambda x) \quad (2.30)$$

$$\mathbf{E} x = \lambda^{-1}, \quad (2.31)$$

$$\mathbf{D} x = \lambda^{-2}. \quad (2.32)$$

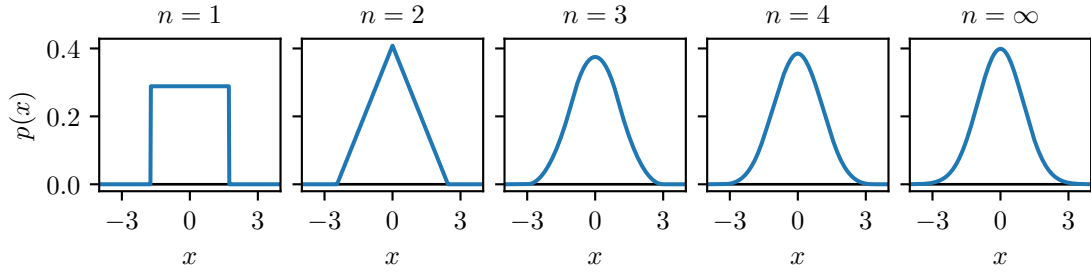
Laplaceova razdioba je kontinuirana razdioba definirana parametrima $\beta \in \mathbb{R}_{>0}$ i $\mu \in \mathbb{R}$ i ima ova svojstva:

$$p(x) = \frac{1}{2\beta} \exp\left(-\frac{|x|}{\beta}\right) \quad (2.33)$$

$$\mathbf{E} x = \mu, \quad (2.34)$$

$$\mathbf{D} x = \beta^2. \quad (2.35)$$

Gaussova (normalna) razdioba je kontinuirana razdioba definirana



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.

parametrima $\mu \in \mathbb{R}$ i $\sigma \in \mathbb{R}_{>0}$ i ima ova svojstva:

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \quad (2.36)$$

$$\mathbf{E} X = \mu, \quad (2.37)$$

$$\mathbf{D} X = \sigma^2. \quad (2.38)$$

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.

$$\lim_{n \rightarrow \infty} P(z_n < z) = \int_{-\infty}^z p_{\mathcal{N}(0,1)}(z') dz'. \quad (2.39)$$

$p_{\mathcal{N}(0,1)}$ označava funkciju gustoće normalne razdiobe s $\mu = 0$ i $\sigma = 1$. To je detaljnije objašnjeno i dokazano npr. u (Elezović, 2007). Centralni granični teorem je ilustriran na slici 2.1.

2.2. Teorija informacije

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

$$I(x \in A) := \log_b \frac{1}{P(x \in A)} = -\log_b P(x \in A). \quad (2.40)$$

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**:

$$H(x) := \mathbf{E}_x I(x = x) = -\mathbf{E}_x \log_b P(x). \quad (2.41)$$

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 (2.41) vjerojatnost zamijeni gustoćom vjerojatnosti, onda on predstavlja **diferencijalnu entropiju**, jedan od analoga² entropije za kontinuirane varijable koji nema neka od svojstava koja ima entropija.

Unakrsna entropija je mjera koja iskazuje donju granicu očekivanja duljine poruke kodirane optimalnim kodom za razdiobu $P(y)$ dok izvor poruka ima razdiobu $P(x)$. Ovako je definirana:

$$H_y(x) := \mathbf{E}_x I(y = x) = -\mathbf{E}_x \log_b P_y(x). \quad (2.42)$$

Za $P(y) = P(x)$ je $H_y(x) = H_x(x) = H(x)$. 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 $H_y(x)$.

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

$$D_{KL}(x \parallel y) := H_y(x) - H(x) = \mathbf{E}_x \log_b \frac{P_x(x)}{P_y(x)}. \quad (2.43)$$

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 $P(y)$.

KL-divergencija će biti 0 akko x i y imaju iste razdiobe. To je ilustrirano slikom 2.2. KL-divergencija, kao ni unakrsna entropija, nije simetrična (slika 2.3), tj. općenito $D_{KL}(x \parallel y) \neq D_{KL}(y \parallel x)$ i $H_y(x) \neq H_x(y)$. KL-divergencija je izrazom (2.43)

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

$H(x)$	$D_{\text{KL}}(x \parallel y)$
$H_y(x)$	

Slika 2.2: Odnos entropije, unakrsne entropije i KL-divergencije.

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:

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

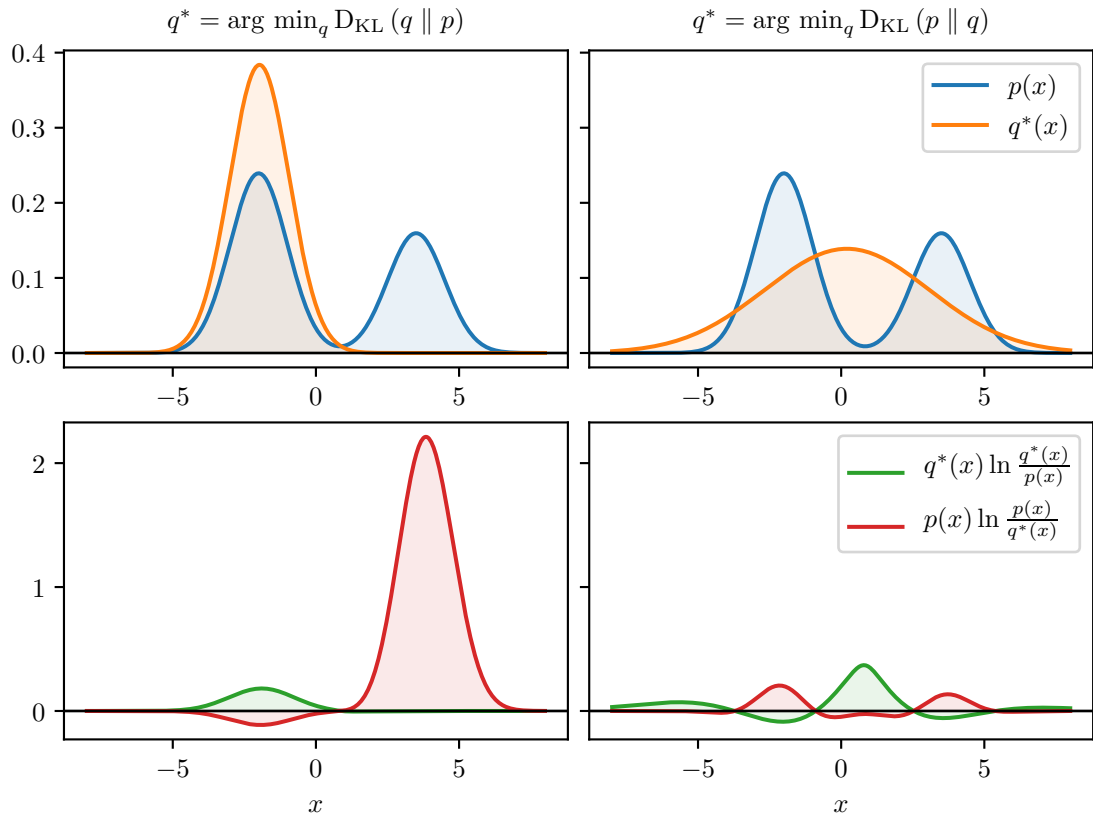
$$= H(x) - H(x | y) \quad (2.46)$$

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

gdje je

$$H(x | y) := \mathbf{E}_x H(y | x = x) \quad (2.48)$$

uvjetna entropija. 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 $I(x; y) = H(x) = H(y)$. Definirane veličine mogu se prikazati kao na slici 2.4. Isti odnosi vrijede ako se entropija zamijeni diferencijalnom entropijom.



Slika 2.3: Asimetričnost KL-divergencije. p je fiksna razdioba (funkcija gustoće), a q^* je Gaussova razdioba koja ju aproksimira minimizacijom KL-divergencije $D_{KL}(q \parallel p)$ (lijevo) ili $D_{KL}(p \parallel q)$ (desno). U donjem retku grafovi prikazuju podintegralne funkcije odgovarajućih KL-divergencija. Kod njih zbrojevi predznačenih površina obojanih područja odgovaraju KL-divergencijama $D_{KL}(q \parallel p)$ (zeleno) ili $D_{KL}(p \parallel q)$ (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).

H(x)		
H(x y)	I(x, y)	H(y x)
	H(y)	
H(x, y)		

Slika 2.4: Odnosi informacijsko-teorijskih veličina dviju slučajnih varijabli.

2.3. Optimizacija temeljena na gradijentu

U ovom odjeljku su opisani osnovni optimizacijski algoritmi temeljeni na gradijentu. Oni su bitni u strojnom učenju (poglavlje 4), posebno u dubokom učenju (poglavlje 5). Izvedeni algoritmi koji se primjenjuju u dubokom učenju opisani su u pododjeljku 5.3.2.

Neka je $f: \mathbb{R}^n \rightarrow \mathbb{R}$ funkcija čiji minumom želimo naći s obzirom na parametre \mathbf{x} . Ona se u okolini točke \mathbf{x} , ako je dovoljno (beskonačno) puta derivabilna može izraziti Taylorovim redom:

$$f(\mathbf{x} + \mathbf{d}) = f(\mathbf{x}) + \frac{\partial}{\partial \mathbf{x}} f(\mathbf{x}) \mathbf{d} + \frac{1}{2} \mathbf{d}^T \frac{\partial^2 f}{\partial \mathbf{x} \partial \mathbf{x}^T} \mathbf{x} f(\mathbf{x}) \mathbf{d} + \dots \quad (2.49)$$

S drugačijim oznakama:

$$f(\mathbf{x} + \mathbf{d}) = f(\mathbf{x}) + \nabla_{\mathbf{x}} f(\mathbf{x})^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \mathbf{H}_f(\mathbf{x}) \mathbf{d} + \dots \quad (2.50)$$

2.3.1. Gradijentni spust

Ako je \mathbf{d} ima malu normu, funkciju f u okoline neke točke možemo dobro aproksimirati s prvih nekoliko članova Taylorovog reda. **Gradijentni spust** je optimizacijski algoritam koji koristi linearnu aproksimaciju i iterativnim ažuriranjem parametara u smjeru gradijenta (*najstrmijem* smjeru) traži minimum. Iteracija gradijentnog spusta ima ovakav oblik:

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \eta \nabla_{\mathbf{x}} f(\mathbf{x}_i), \quad (2.51)$$

gdje je i redni broj iteracije, a η **veličina koraka** (**stopa učenja** kod strojnog učenja) koja može biti konstanta ili može ovisiti o broju iteracije i . Neka $\mathbf{g} = \nabla_{\mathbf{x}} f(\mathbf{x})$ i $\mathbf{H} = \mathbf{H}_f(\mathbf{x})$. Za dovoljno mal η

$$f(\mathbf{x} - \eta \mathbf{g}) \approx f(\mathbf{x}) - \eta \mathbf{g}^T \mathbf{g} - \frac{1}{2} \eta^2 \mathbf{g}^T \mathbf{H} \mathbf{g} \quad (2.52)$$

Uz neke blage uvjete koje mora zadovoljavati f i dovoljno mal η , gradijentni spust konvergira, tj. proizvoljno se blizu približi nekom lokalnom minimumu (ili stacionarnoj točki koja nije lokalni minimum, gdje $\nabla_{\mathbf{x}} f(\mathbf{x}) = \mathbf{0}$) ovisno o η . Jedan blagi uvjet može biti **Lipschitz kontinuiranost** funkcije f ili njene derivacije (Goodfellow et al., 2016). Funkcija f je Lipschitz kontinuirana ako postoji

konstanta λ za koju za svaki par (\mathbf{x}, \mathbf{y}) vrijedi:

$$|f(\mathbf{x}) - f(\mathbf{y})| < \lambda \|\mathbf{x} - \mathbf{y}\|. \quad (2.53)$$

Najmanji takav λ naziva se **Lipschitzova konstanta**.

2.3.2. Postupci drugog reda

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

Ako koristimo kvadratnu aproksimaciju (2.52), možemo pokušati pronaći optimalni η koji ju minimizira. η za koji $\frac{\partial}{\partial \eta} f(\mathbf{x} - \eta \mathbf{g}) = 0$ će, ako $\mathbf{g}^\top \mathbf{H} \mathbf{g} > 0$ dati minimum u smjeru gradijenta kvadratne aproksimacije funkcije f u točki \mathbf{x} . Dobije se:

$$\eta = \frac{\mathbf{g}^\top \mathbf{g}}{\mathbf{g}^\top \mathbf{H} \mathbf{g}}. \quad (2.54)$$

Ako je $f: \mathbb{R}^n \rightarrow \mathbb{R}$ konveksna (pozitivno definitna) kvadratna funkcija, izmijenjeni algoritam gradijentnog spusta, koji ovako određuje veličinu koraka, minimum pronalazi u najviše n koraka.

Postupak drugog reda koji se ne ograničava na pomake u smjeru gradijenta je **Newton-Raphsonov postupak**. Deriviranjem desne strane jednadžbe (2.50) po \mathbf{d} i izjednačavanjem s 0 dobiva se:

$$\mathbf{0} = \nabla_{\mathbf{x}} f(\mathbf{x})^\top + \mathbf{d}^\top \mathbf{H}_f(\mathbf{x}) + \dots. \quad (2.55)$$

Uz kvadratnu aproksimaciju i kraće oznake $\mathbf{g} = \nabla_{\mathbf{x}} f(\mathbf{x})$ i $\mathbf{H} = \mathbf{H}_f(\mathbf{x})$: $\mathbf{H} \mathbf{d} = -\mathbf{g}$. Slijedi da je pomak \mathbf{d} koji daje stacionarnu točku aproksimacije

$$\mathbf{d} = -\mathbf{H}^{-1} \mathbf{g}. \quad (2.56)$$

Za nekvadratne funkcije, koje imaju pozitivno definitnu Hesseovu matricu u svakoj točki, može se iterativno primjenjivati

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \eta \mathbf{H}_f(\mathbf{x}_i) \nabla_{\mathbf{x}} f(\mathbf{x}_i) \quad (2.57)$$

s $\eta < 1$.

3. Statističko modeliranje

3.1. Probabilistički grafički modeli

Neka su x_1, \dots, x_n slučajne varijable čiju združenu razdiobu razmatramo. Želimo na temelju opežanih varijabli korištenjem pravila vjerojatnosti **zaključivati** o razdiobama nekih neopažanih varijabli. Općenito, zaključivanje se provodi uvjetovanjem po opažanim varijablama i marginalizacijom po varijablama koje nas ne zanimaju izravno (Murphy, 2012):

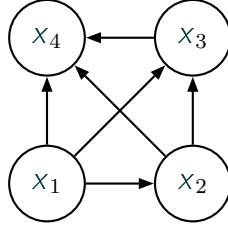
$$p(\mathbf{x}_q | \mathbf{x}_o) = \frac{p(\mathbf{x}_q, \mathbf{x}_o)}{p(\mathbf{x}_o)} = \frac{\int p(\mathbf{x}_q, \mathbf{x}_n, \mathbf{x}_o) d\mathbf{x}_n}{\int p(\mathbf{x}_q, \mathbf{x}_n, \mathbf{x}_o) d(\mathbf{x}_q, \mathbf{x}_n)}. \quad (3.1)$$

Ovdje je \mathbf{x}_q niz varijabli o kojima želimo zaključivati (varijable upita), \mathbf{x}_o niz opažanih varijabli, a \mathbf{x}_n niz varijabli *smetnje* (*nuisance*).

Zavisnosti između slučajnih varijabli otežavaju modeliranje i zaključivanje – potrebno je više podataka i zaključivanje je računski zahtjevnije. Obično možemo pretpostaviti uvjetne zavisnosti između slučajnih varijabli, što se može predstaviti neusmjerenim ili usmjerenim grafom. Prema definiciji na Wikipediji ¹, **probabilistički grafički model** ili **grafički model** je probabilistički model koji se može prikazati grafom koji izražava strukturu uvjetnih zavisnosti među slučajnim varijablama. U tom grafu čvorovi označavaju slučajne varijable, a bridovi zavisnosti. Usmjereni bridovi označavaju modeliranje uvjetne zavisnosti, a neusmjereni združeno modeliranje. Ako je graf grafičkog modela usmjeren i acikličan, on se naziva **Bayesova mreža** ili **Bayesovski model**, a ako je neusmjeren, naziva se **Markovljeva mreža** ili **Markovljevo slučajno polje** (engl. *Markov random field*, *MRF*). U nastavku ovog odjeljka naglasak će biti na Bayesovim mrežama.

Združena razdioba se prema pravilu umnoška (jednadžba 2.6) može npr. ovako

¹https://en.wikipedia.org/wiki/Graphical_model



Slika 3.1: Prikaz grafičkog modela s faktorizacijom
 $p(x_1, x_2, x_3, x_4) = p(x_1) p(x_2 | x_1) p(x_3 | x_1, x_2) p(x_4 | x_1, x_2, x_3)$.

izraziti:

$$p(x_1, \dots, x_n) = p(x_1) p(x_2 | x_1) \cdots p(x_n | x_1, \dots, x_{n-1}) \quad (3.2)$$

$$= \prod_i p(x_i | x_1, \dots, x_{i-1}). \quad (3.3)$$

Prema tome, svaki probabilistički grafički model ima ekvivalentnu Bayesovu mrežu. Ako uzmemo $n = 4$, graf koji odgovara faktorizaciji u jednadžbi (3.5) prikazan je na slici 3.1.

Pretpostavljanjem uvjetnih nezavisnosti, neki bridovi grafa G se mogu ukloniti pa za varijable (čvorove grafa) vrijedi **uređajno Markovljevo svojstvo**:

$$x \perp \text{pred}_G(x) \setminus \text{pa}_G(x) \mid \text{pa}_G(x). \quad (3.4)$$

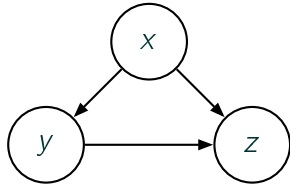
Jednadžba (3.5) onda prelazi u

$$p(x_1, \dots, x_n) = \prod_i p\left(x_i \mid \bigcap_{x_j \in \text{pa}_G(x_i)} \{x_j = x_j\}\right). \quad (3.5)$$

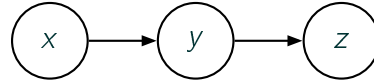
To omogućuje primjenu efikasnijih algoritama za zaključivanje (Murphy, 2012). Na slici 3.2 prikazani su osnovni slučajevi odnosa između triju slučajnih varijabli povezanih zavisnostima koje mogu biti dio većeg grafa. Oni su detaljnije objašnjeni npr. u Bishop (2006) i Alpaydin (2014).

Na slici 3.3 prikazan je primjer na kojemu se koriste još neke oznake: sivi čvorovi označavaju opažane varijable, četverokut označava veći broj podgrafova s istom strukturom.

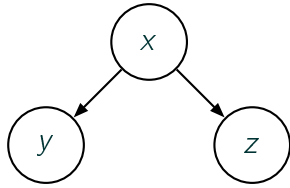
Općenitije, o uvjetnoj nezavisnosti podskupova varijabli govori svojstvo **d-separacije**. Kažemo da je staza (podgraf sa strukturom lanca) P grafa G **d-odvojena** skupom čvorova \mathbb{E} akko P sadrži barem jedno od sljedećeg (Murphy, 2012):



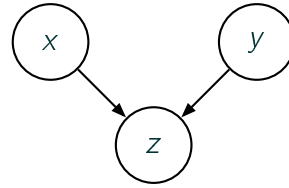
(a) Grafički model s faktorizacijom $p(x, y, z) = p(x) p(y | x) p(z | x, y)$.



(b) Uz $x \perp z | y$ faktorizacija postaje $p(x, y, z) = p(x) p(y | x) p(z | y)$ (lanac).

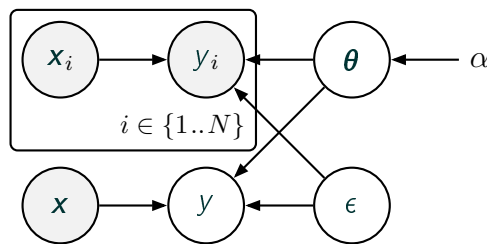


(c) Uz $y \perp z | x$ faktorizacija postaje $p(x, y, z) = p(x) p(y | x) p(z | x)$ (račvanje).



(d) Uz $x \perp y$ faktorizacija postaje $p(x, y, z) = p(x) p(y) p(z | x, y)$ (sraz). Ovdje također vrijedi $x \not\perp y | z$.

Slika 3.2: Osnovni slučajevi uvjetne nezavisnosti. Slike b, c i d prikazuju grafove dobivene uvođenjem pretpostavki uvjetne nezavisnosti za grafički model s 3 slučajne varijable prikazan na slici a.



Slika 3.3: Primjer grafičkog modela s faktorizacijom $p(x, y, x_1 \dots x_N, y_1 \dots y_N, \theta, \epsilon) = p(\theta) p(\epsilon) p_x(x) p_{y|x, \theta, \epsilon}(y | x, \theta, \epsilon) \prod_i (p_x(x_i) p_{y_i|x_i, \theta, \epsilon}(y_i | x_i, \theta, \epsilon))$. Graf prikazuje model regresije, gdje su θ nepoznati parametri, x_i i y_i opažani parovi ulaza i izlaza, x opažani ulaz s nepoznati izlazom, a ϵ homoskedastički šum, tj. šum koji ne ovisi o ulazu. Na slici je još eksplicitno prikazana deterministička varijabla α koja je parametar razdiobe $p(\theta) = p(\theta | \alpha)$. Slika je napravljena po uzoru na sliku 14.7 u [Alpaydin \(2014\)](#).

- lanac $a \rightarrow b \rightarrow c$, gdje $b \in E$
- račvanje $a \leftarrow b \rightarrow c$, gdje $b \in E$
- sraz $a \rightarrow b \leftarrow c$, gdje $\forall b' \in \{b\} \cup \text{succ}_G(b), b' \notin E$.

Kažemo da skup čvorova E d-odvaja čvorove x i y akko su sve staze između njih d-odvojene. Vrijedi $x \perp y \mid E$ akko skup čvorova E d-odvaja čvorove x i y . To se može poopćiti na skupove čvorova. Skup čvorova opažanjem kojega neki čvor postaje neovisan o ostatku grafa naziva se **Markovljev pokrivač** (engl. *Markov blanket*). Markovljev pokrivač čvora x je

$$\text{pa}_G(x) \cup \text{ch}_G(x) \cup \bigcup_{y \in \text{ch}_G(x)} \text{pa}_G(y) \quad (3.6)$$

U navedenim knjigama opisani su algoritmi koji se koriste za efikasno zaključivanje iskorištavanjem strukture grafa.

3.2. Procjena parametara i zaključivanje

3.2.1. Procjenitelji i točkaste procjene parametara

Ovaj pododjeljak se temelji na [Elezović \(2007\)](#).

Neka je x slučajna varijabla i $p(x)$ 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

$$\hat{\theta} = f(x_1, \dots, x_N). \quad (3.7)$$

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

$$\hat{\theta} = f(\mathcal{D}). \quad (3.8)$$

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

3.2.2. Svojstva i pogreška procjenitelja

Priistranost 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

$$\mathbf{E} \hat{\theta} = \theta. \quad (3.9)$$

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

$$\lim_{N \rightarrow \infty} \mathbf{D} [\hat{\theta}(\mathcal{D}_N)] = 0. \quad (3.10)$$

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

$$\mathbf{E} [(\hat{\theta} - \theta)^2] = \mathbf{D} \hat{\theta} + (\mathbf{E} \hat{\theta} - \theta)^2. \quad (3.11)$$

3.2.3. Procjenitelj maksimalne izglednosti

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

$$\theta_{\text{ML}} = \arg \max_{\theta} p(\mathcal{D} | \theta). \quad (3.12)$$

Zbog pretpostavke međusobne nezavisnosti primjera vrijedi

$$p(\mathcal{D} | \theta) = \prod_{d \in \mathcal{D}} p(d | \theta). \quad (3.13)$$

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

$$p(\mathcal{D} | \theta) = \prod_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} p(\mathbf{y} | \mathbf{x}, \theta) p(\mathbf{x} | \theta) = p(\mathbf{x}) \prod_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} p(\mathbf{y} | \mathbf{x}, \theta). \quad (3.14)$$

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

3.2.4. Procjenitelj maksimalne aposteriorne vjerojatnosti

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

$$p(\theta | \mathcal{D}) = \frac{p(\mathcal{D} | \theta) p(\theta)}{p(\mathcal{D})} = \frac{p(\mathcal{D} | \theta) p(\theta)}{\int p(\mathcal{D} | \theta') p(\theta') d\theta'}. \quad (3.15)$$

Maksimizacijom aposteriorne vjerojatnosti dobivaju se parametri

$$\theta_{\text{MAP}} = \arg \max_{\theta} p(\theta | \mathcal{D}) = \arg \max_{\theta} p(\mathcal{D} | \theta) p(\theta). \quad (3.16)$$

Ovdje nije potrebno normalizirati aposteriornu vjerojatnost izračunavanjem **marginalne izglednosti** (engl. *marginal likelihood, evidence*) $p(\mathcal{D})$ u nazivniku na desnoj strani jednadžbe (3.15) jer ona ne ovisi θ , nego samo o modelu \mathcal{H} . Odabirom uniformne (neinformativne) apriorne razdiobe MAP-procjenitelj postaje ekvivalentan ML-procjenitelju.

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

3.2.5. 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 $p(\theta | \mathcal{D})$ za koju je integriranjem po svim mogućim parametrima potrebno izračunati marginalnu izglednost $p(\mathcal{D}) = \int p(\mathcal{D} | \theta') p(\theta') d\theta'$ iz nazivnika na desnoj strani jednadžbe (3.23).

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 parametrima (Neal, 1995):

$$p(\mathbf{d} | \mathcal{D}) = \int p(\mathbf{d} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathcal{D}) d\boldsymbol{\theta} = \mathbf{E}_{\boldsymbol{\theta} | \mathcal{D}} p(\mathbf{d} | \boldsymbol{\theta}), \quad (3.17)$$

gdje je korištena uvjetna nezavisnost $\mathbf{d} \perp \mathcal{D} | \boldsymbol{\theta}$.

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

$$p(\mathbf{d} | \mathcal{D}) \approx \int p(\mathbf{d} | \boldsymbol{\theta}) \delta(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) d\boldsymbol{\theta} = p(\mathbf{d} | \hat{\boldsymbol{\theta}}). \quad (3.18)$$

Za diskriminativne modele se bayesovsko zaključivanje može izraziti ovako:

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

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

$$p(\mathbf{y} | \mathbf{x}, \mathcal{D}) = \int p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathcal{D}) d\boldsymbol{\theta} = \mathbf{E}_{\boldsymbol{\theta} | \mathcal{D}} p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}). \quad (3.19)$$

Kod regresije je često, ako pretpostavljamo da pogreška izlaza ima Gaussovu razdiobu, najbolja procjena hipoteze očekivanje po naučenoj razdiobi parametara (Neal, 1995):

$$h(\mathbf{x}) = \mathbf{E}_{\boldsymbol{\theta} | \mathcal{D}} h(\mathbf{x}; \boldsymbol{\theta}) = \int h(\mathbf{x}; \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathcal{D}) d\boldsymbol{\theta}. \quad (3.20)$$

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

3.3. Monte Carlo aproksimacija

Ovaj odjeljak se temelji 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 integral (ili zbroj). Ako se on može ovako izraziti:

$$s = \int p(x) f(x) dx = \mathbf{E} f(x), \quad (3.21)$$

može se procijeniti uzorkovanjem:

$$\hat{s}_n = \frac{1}{n} \sum_{i=1}^n f(x_i). \quad (3.22)$$

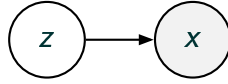
Procjenitelj \hat{s}_n je nepristran ako su x_i nezavisne i imaju istu razdiobu kao x i valjan ako su varijance $f(x_i)$ ograničene. Vrijedi $\mathbf{D} \hat{s}_n = \frac{1}{n} \mathbf{D} f(x)$.

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

3.4. Aproksimacija razdioba i aproksimacijsko zaključivanje

Ovaj odjeljak se uglavnom 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 modela aposteriorna razdioba se ne može lako izračunati i treba koristiti aproksimacijske postupke od kojih su glavni **varijacijski** postupci ([Jordan et al., 1999](#)) i postupci **Monte Carlo** aproksimacije s uzorkovanjem pomoću **Markovljevog lanca** (MCMC, engl. *Markov chain Monte Carlo*). MCMC-postupci temelje se na definiranju stohastičkog procesa koji ima stacionarnu razdiobu jednaku razdiobi koja se aproksimira, omogućuju asimptotski egzaktno uzorkovanje velike klase razdioba. Varijacijski postupci temelje se na aproksimaciji razdiobe nekom jednostavnijom koja se pronalazi rješavanjem optimizacijskog problema, brži su i jednostavniji za ostvariti za složenije modele.



Slika 3.4: Prikaz grafičkog modela sa skrivenom varijablom z i opažanom varijablom x .

Razmatramo bayesovski model koji ima latentnu varijablu z i vidljivu varijablu x . Model je prikazan na slici 3.4 i opisan je ovom jednažbom združene vjerojatnosti:

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

Zaključivanjem se određuje aposteriorna razdioba latentne varijable:

$$p(\mathbf{z} | \mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})} = \frac{p(\mathbf{x}, \mathbf{z})}{\int p(\mathbf{x}, \mathbf{z}) d\mathbf{z}}. \quad (3.23)$$

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

3.5. Postupci uzorkovanja

TODO

1.3.1 Neal (1995).

3.6. Varijacijsko zaključivanje

Za razliku od uzorkovanja kod MCMC-postupaka, osnovna ideja kod varijacijskog zaključivanja je optimizacija. Prvo se odabire familija razdioba

$\mathcal{Q} = \{p(\tilde{\mathbf{z}})\}_{\tilde{\mathbf{z}}} = \{q_{\phi}\}_{\phi}$ 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 ϕ :

$$q^* = \arg \min_{p(\tilde{\mathbf{z}}) \in \mathcal{Q}} D_{\text{KL}}(\tilde{\mathbf{z}} \parallel (\mathbf{z} | \mathbf{x})). \quad (3.24)$$

Naziv **varijacijsko zaključivanje** dolazi od varijacijskog računa², gdje se koriste varijacije, tj. male promjene u funkcijama i funkcionalima, kako bi se pronašli minimumi ili maksimumi funkcionala, preslikavanja iz skupa funkcija u \mathbb{R} , koji su često izraženi kao integrali koji uključuju funkcije i njihove derivacije.

Neka je q oznaka funkcije gustoće vjerojatnosti aproksimacijske razdiobe: $q := p_{\tilde{z}}$. Ako ciljnu funkciju ovako izrazimo:

$$\begin{aligned} D_{\text{KL}}(\tilde{z} \parallel (z \mid \mathbf{x})) &= \mathbf{E}_{\tilde{z}} \ln \frac{q(\tilde{z})}{p_{z \mid \mathbf{x}}(\tilde{z})} \\ &= \mathbf{E}_{\tilde{z}} \ln q(\tilde{z}) - \mathbf{E}_{\tilde{z}} \ln p(z = \tilde{z}, \mathbf{x}) + \ln p(\mathbf{x}), \end{aligned} \quad (3.25)$$

vidi se da se ona ne može lako izračunati jer zahtijeva računanje marginalne izglednosti $p(\mathbf{x})$ iz nazivnika u jednadžbi (3.23) marginalizacijom po z . Marginalna izglednost ne ovisi o varijacijskim parametrima pa ju možemo zanemariti i maksimiziramo funkciju koja se naziva **varijacijska donja granica** (engl. *variational lower bound*) ili **donja granica (logaritma) marginalne izglednosti** (engl. *(log) evidence lower bound, ELBO*):

$$L_{\mathbf{x}}(\tilde{z}) := \ln p(\mathbf{x}) - D_{\text{KL}}(\tilde{z} \parallel (z \mid \mathbf{x})) = \mathbf{E}_{\tilde{z}} \ln p(z = \tilde{z}, \mathbf{x}) - \mathbf{E}_{\tilde{z}} \ln q(\tilde{z}). \quad (3.26)$$

Ona se može i ovako izraziti:

$$L_{\mathbf{x}}(\tilde{z}) = \mathbf{E}_{\tilde{z}} \ln p(\mathbf{x} \mid z = \tilde{z}) - D_{\text{KL}}(\tilde{z} \parallel z). \quad (3.27)$$

Maksimiziranje takve ciljne funkcije s obzirom na varijacijske parametre daje razdiobu $q^* = p(\tilde{z}^*)$ koja dobro objašnjava podatke jer se potiče veće očekivanje logaritma izglednosti (prvi član), a 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 $\ln p(\mathbf{x}) \geq L_{\mathbf{x}}(\tilde{z})$ preko Jensenove nejednakosti. Ta nejednakost slijedi i iz prethodne jednadžbe i nenegativnosti KL-divergencije:

$$\ln p(\mathbf{x}) = L_{\mathbf{x}}(\tilde{z}) + D_{\text{KL}}(\tilde{z} \parallel (z \mid \mathbf{x})) \geq L_{\mathbf{x}}(\tilde{z}). \quad (3.28)$$

²https://en.wikipedia.org/wiki/Calculus_of_variations

3.6.1. Metoda polja sredina

Dodatno pojednostavljenje koje pomaže u modeliranju i optimizaciji je pretpostavljanje nezavisnosti između latentnih varijabli. Onda za varijacijsku razdiobu vrijedi ovakva faktorizacija:

$$q(\tilde{\mathbf{z}}) = \prod_i q_i(\tilde{z}_i), \quad (3.29)$$

gdje su q_i funkcije gustoće pojedinih slučajnih varijabli, a $\tilde{z}_i = \tilde{\mathbf{z}}_{[i]}$. Kod **metode polja sredina** pretpostavlja se takva razdioba i obično se primjenjuje koordinatni spust za optimizaciju, s čime ima veze ime metode. To je detaljnije objašnjeno u [Murphy \(2012\)](#).

4. Nadzirano strojno učenje

Ovo poglavlje se uglavnom temelji na Šnajder i Dalbelo Bašić (2014) i Goodfellow et al. (2016).

Zadatak algoritama **nadziranog učenja** je preslikavanje **ulaznih primjera** x iz **ulaznog prostora** \mathbb{X} u **izlaze (oznake)** $y \in \mathbb{Y}$ na temelju konačnog skupa označenih primjera $\mathcal{D} = \{(x_i, y_i)\}_i$. Algoritmima strojnog učenja pretražuje se **model** ili **prostor hipoteza** u cilju pronalaska **hipoteze** koja osim primjera iz skupa za učenje, u izlaze dobro preslikava i primjere koji nisu u skupu za učenje. Sposobnost postizanja dobre performanse na neviđenim primjerima naziva se **generalizacija**.

Neka je $\mathcal{D} = \{d_i\}_i$ 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: $p(d \mid \mathcal{D}, \mathcal{H})$. Model koji modelira razdiobu primjera nazivamo **generativnim modelom**. U nastavku ćemo izostavljati oznaku modela radi kraćeg zapisa.

Ako su primjeri parovi $d_i = (x_i, y_i) \in \mathbb{X} \times \mathbb{Y}$, 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 $p(y \mid x)$ hipotezom koja ulazni primjer x preslikava u razdiobu $p(y \mid x, \mathcal{D})$. Neprobabilistički diskriminativni modeli modeliraju funkciju dodjeljivanja oznaka hipotezom $h: \mathbb{X} \rightarrow \mathbb{Y}$. Modeliranje zajedničke razdiobe $p(x, y)$ obično zahtijeva više računalnih resursa i podataka (Bishop, 2006).

Modeli se još mogu podijeliti na **parametarske** i **neparametarske**. Kod parametarskih modela broj parametara je unaprijed određen, dok kod neparametarskih on ovisi o podacima za učenje.

4.1. Induktivna pristranost

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.

4.2. Komponente algoritma strojnog učenja

Prema Šnajder i Dalbelo Bašić (2014), kod većine algoritama strojnog učenja možemo razlikovati 3 osnovne komponente, 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 **funkciju 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:

$$E(\theta, \mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} L(\mathbf{y}, h(\mathbf{x}; \theta)). \quad (4.1)$$

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 4.5.

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

$$\theta^* = \arg \min_{\theta} E(\theta, \mathcal{D}). \quad (4.2)$$

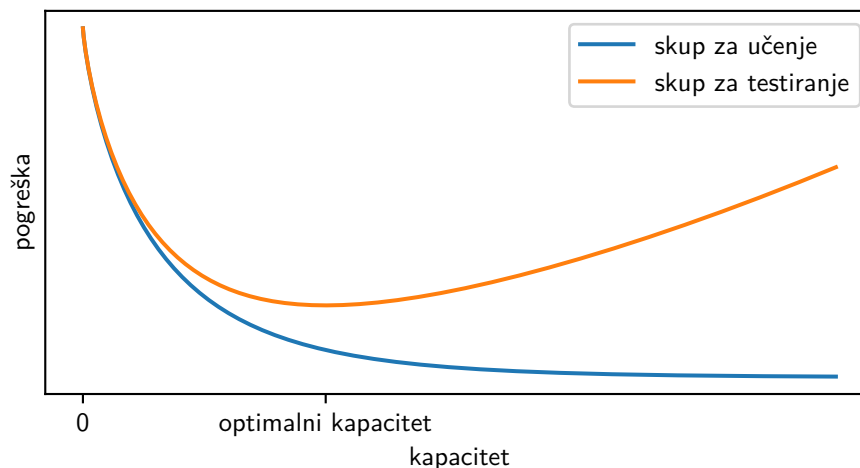
Kod nekih jednostavnijih modela minimum možemo odrediti analitički. Inače moramo koristiti neki iterativni optimizacijski postupak. Kod nekih složenijih modela, kao što su neuronske mreže, funkcija pogreške nije unimodalna i vjerojatnost pronalaska globalnog optimuma je zanemariva, ali ipak se mogu pronaći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

U literaturi riječ *model* često ima šire značenje. Uz skup hipoteza obično obuhvaća i induktivnu pristranost ili dio nje. Model u tom smislu bi se formalno mogao definirati kao par (\mathcal{H}, B) , gdje je \mathcal{H} skup mogućih hipoteza, a B induktivna pristranost koja hipotezama dodjeljuje različite važnosti. Ovdje će se u nastavku koristiti takvo značenje riječi *model*, a riječ *prostor hipoteza* će se koristiti sa značenjem modela u užem smislu.

4.3. Kapacitet modela, podnaučenost i prenaučenost

Cilj algoritama strojnog učenja je postići malu **pogrešku generalizacije**, tj. malo očekivanje pogreške na primjera koji nisu korišteni za učenje i odabir modela. Generalizacijska pogreška se procjenjuje pogreškom na **skupu za testiranje**. Obično pretpostavljamo da su skupovi primjera koje koristimo za učenje, odabir modela i testiranje generirani međusobno nezavisno i iz iste razdiobe.

Kapacitet ili složenost modela je svojstvo koje opisuje njegovu sposobnost prilagodbe podacima. Model koji se previše prilagođava podacima za učenje (i statističkom šumu u njima) obično ima slabu prediktivnu moć. Treba odabrati model (ili hipotezu) koji dobro objašnjava podatke, ali nije previše složen. O tome govori i načelo **Occamove oštrice** prema kojemu među hipotezama konzistentnima s opažanjem treba odbaciti sve osim najjednostavnije od njih. Postoje formalizacije Occamove oštrice (Blumer et al., 1987, 1989; Grünwald, 2005; Rathmanner i



Slika 4.1: Ovisnost pogrešaka na skupovima za učenje i testiranje o kapacitetu modela. Povećavanjem kapaciteta povećava se razlika između pogreške na skupovima za testiranje i pogreške na skupovima za učenje.

Hutter, 2011). Na ograničavanje složenosti modela možemo utjecati ograničavanjem prostora hipoteza i regularizacijom (*mekim* ograničavanjem).

Model s većim kapacitetom (složeniji model) može postići manju pogrešku na skupovima za učenje. Prevelik kapacitet povećava pogrešku generalizacije. Za model koji daje veliku pogrešku generalizacije kažemo da je **prenaučen**. Kod takvog modela hipoteze će jako varirati u ovisnosti o skupovima za učenje i zato kažemo da složeni modeli imaju visoku varijancu. Model premalog kapaciteta (prejednostavan model) ima manju razliku između pogreške na skupovima za učenje i pogreške na skupovima za testiranje, ali su obje pogreške veće od optimalnih. Za model koji ne postiže malu pogrešku na skupovima za učenje kažemo da je **podnaučen**. U jednostavan model ugrađene su jače pretpostavke i kažemo da on ima jaču pristranost. Uobičajena ovisnost pogrešaka na skupovima za učenje i testiranje o kapacitetu ilustrirana je slikom 4.1.

4.4. Odabir modela

TODO

Murray i Ghahramani (2005) MacKay

4.5. Funkcija pogreške

Dijelovi ovog odjeljka temelje se na (Murphy, 2012).

4.5.1. 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. **Rizik** definiramo kao očekivanje funkcije gubitka:

$$R(\theta; \mathcal{D}) = \mathbf{E}_{(x,y) \sim \mathcal{D}} L(y, h(x; \theta)). \quad (4.3)$$

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

$$R_E(\theta; \mathbb{D}) = \mathbf{E}_{(x,y) \sim \mathbb{D}} L(y, h(x; \theta)) = \frac{1}{|\mathbb{D}|} \sum_{(x,y) \in \mathbb{D}} L(y, h(x; \theta)). \quad (4.4)$$

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

$$R(\theta; \mathcal{D}) = \mathbf{E}_{d \sim \mathcal{D}} L(d, h(d; \theta)). \quad (4.5)$$

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

$$R_E(\theta; \mathbb{D}) = -\frac{1}{|\mathbb{D}|} \ln p(\mathbb{D} | \theta) = -\frac{1}{|\mathbb{D}|} \sum_{d \in \mathbb{D}} \ln p(d | \theta), \quad (4.6)$$

gdje je korištena pretpostavka međusobne nezavisnosti primjera. Gubitak je onda $L(d, h(d; \theta)) = -\ln p(d | \theta)$. U slučaju diskriminativnog modela, uz zanemarivanja faktora izglednosti koji ne ovisi o θ (jednadžba (3.14)), vrijedi $L(d, h(x; \theta)) = -\ln p(y | x, \theta)$. Minimizacija gubitka definiranog kao negativni logaritam izglednosti ekvivalentna je minimizaciji KL-divergencije ili unakrsne entropije (odjeljak 2.2) s obzirom na empirijsku razdiobu. Zbog toga se takav gubitak još naziva **gubitak unakrsne entropije**.

4.5.2. Strukturni rizik i regularizacija

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** $\lambda R_R(\theta)$, $\lambda \geq 0$, koji predstavlja **strukturni rizik** koji daje prednost jednostavnijim hipotezama. Funkcija pogreške onda ima ovakav oblik:

$$E(\theta; \mathcal{D}) = R_E(\theta; \mathcal{D}) + \lambda R_R(\theta). \quad (4.7)$$

Regularizacijski gubitak obično ovisi samo o parametrima, ali može ovisiti i o podacima (Goodfellow et al., 2016).

Ako kao funkciju pogreške koristimo negativni logaritam aposteriorne vjerojatnosti uz pretpostavku međusobne nezavisnosti primjera, funkcija pogreške je

$$E(\theta; \mathcal{D}) = -\frac{1}{|\mathcal{D}|} \ln p(\theta | \mathcal{D}) \quad (4.8)$$

$$= \underbrace{-\frac{1}{|\mathcal{D}|} \ln p(\mathcal{D} | \theta)}_{R_E(\theta; \mathcal{D})} - \underbrace{\frac{1}{|\mathcal{D}|} \ln p(\theta)}_{\lambda R_R(\theta)} + C_1, \quad (4.9)$$

gdje je $C_1 = \frac{1}{|\mathcal{D}|} \ln p(\mathcal{D})$ konstanta koja ne ovisi o θ . Hiperparametar λ je onda parametar apriorne razdiobe. Možemo ga ovako izlučiti:

$$\ln p(\theta) = \lambda \ln p_0(\theta) + C_2 = \ln p_0(\theta)^\lambda + C_2, \quad (4.10)$$

gdje je $C_2 = -\ln\left(\int_{\theta'} p_0(\theta') d\theta'\right)$ konstanta koja ne ovisi o θ . Vidi se da λ određuje koncentraciju apriorne razdiobe. Povećanje λ smanjuje entropiju apriorne razdiobe. Ona postaje koncentriranija i regularizacija jača. Jačom regularizacijom se povećava pristranost i smanjuje varijanca.

4.6. Osnovni zadaci nadziranog učenja

Osnovni zadaci nadziranog učenja su **klasifikacija** i **regresija**. Zadatak klasifikacije je svakom ulaznom primjerima dodjeljivati oznake iz konačnog skupa oznaka, npr. $\{1..C\}$, gdje svaka oznaka predstavlja jednu **klasu (razred)**. Zadatak regresije je ulaznim primjerima dodjeljivati vrijednosti iz kontinuiranog skupa (obično

\mathbb{R} ili \mathbb{R}^n). Ulazni primjeri su obično realni vektori. Klasifikacijska hipoteza se može definirati preko funkcije s kontinuiranom kodomenom. Ako $C = 2$, ta funkcija može biti $h: \mathbb{X} \rightarrow \mathbb{R}$, a hipoteza $h_c(\mathbf{x}) = \llbracket h(\mathbf{x}) > 0 \rrbracket$. Ako $C > 2$, onda to može biti npr. $h_c(\mathbf{x}) = \arg \max_i h_i(\mathbf{x})$, gdje $h: \mathbb{X} \rightarrow \mathbb{R}^C$ i $h(\mathbf{x}) = [h_i(\mathbf{x})]_{i=1..C}^T$. Kod nekih zadataka ulazi ili izlazi imaju složeniju strukturu i ona se može razlikovati između različitih primjera.

4.7. Primjeri modela: poopćeni linearni modeli

Ovaj odjeljak se temelji na (Šnajder i Dalbello Bašić, 2014).

Linearni modeli su modeli kod kojih je hipoteza definirana afinom transformacijom:

$$h(\mathbf{x}) = h(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{w}^T \mathbf{x} + b, \quad (4.11)$$

gdje je \mathbf{w} vektor **težina**, b **pomak** (engl *bias*), a $\boldsymbol{\theta} = (\mathbf{w}, b)$. Kod linearnih modela je, u slučaju klasifikacije, granica $(n - 1)$ -dimenzionalna hiperravnina s normalom \mathbf{w} . Obično se na ulazne primjere primjenjuje neka nelinearna transformacija

$$\begin{aligned} \phi: \mathbb{R}^n &\rightarrow \mathbb{R}^m \\ \mathbf{x} &\mapsto [\phi_1(\mathbf{x}), \dots, \phi_m(\mathbf{x})]^T \end{aligned}$$

koja predstavlja preslikavanje ulaznog prostora u **prostor značajki**. Funkcije $\phi_i: \mathbb{R}^n \rightarrow \mathbb{R}$ nazivaju se **bazne funkcije**. Hipoteza linearnog modela onda ima oblik

$$h(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}). \quad (4.12)$$

Ovdje je izostavljen pomak b jer jedan od izlaza transformacije ϕ može biti konstanta 1 koja se množi s jednom težinom iz \mathbf{w} .

Poopćeni linearni modeli su modeli kod kojih je hipoteza ovako definirana:

$$h(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x})). \quad (4.13)$$

U odnosu na linearne modele, oni još imaju **prijenosnu (aktivacijsku)** funkciju $f: \mathbb{R} \rightarrow \mathbb{R}$. Ako je f nelinearna, model je nelinearan u parametrima, ali granica klasifikacijskog modela je i dalje hiperravnina.

Slijedi pregled nekih linearnih modela prema Šnajder (2017) uz oznake

$s = \mathbf{w}^\top \phi(\mathbf{x})$ i $\mathbf{s} = \mathbf{W} \phi(\mathbf{x})$:

– Linearna regresija:

$$\begin{aligned} h(\mathbf{x}; \mathbf{w}) &= f(s) = s, \\ p(y | \mathbf{x}, \mathbf{w}) &= \mathcal{N}(h(\mathbf{x}), \sigma^2)(y), \\ L(y, h(\mathbf{x})) &= (h(\mathbf{x}) - y)^2, \\ \nabla_{\mathbf{w}} L(y, h(\mathbf{x})) &= (h(\mathbf{x}) - y) \phi(\mathbf{x}), \end{aligned}$$

gdje $y \in \mathbb{R}$.

– Logistička regresija:

$$\begin{aligned} h(\mathbf{x}; \mathbf{w}) &= f(s) = \frac{1}{1 + \exp(s)} = P(y = 1 | \mathbf{x}, \mathbf{w}), \\ P(y | \mathbf{x}, \mathbf{w}) &= h(\mathbf{x})^y (1 - h(\mathbf{x}))^{1-y}, \\ L(y, h(\mathbf{x})) &= -y \ln h(\mathbf{x}) - (1 - y) \ln(1 - h(\mathbf{x})), \\ \nabla_{\mathbf{w}} L(y, h(\mathbf{x})) &= (h(\mathbf{x}) - y) \phi(\mathbf{x}), \end{aligned}$$

gdje $y \in 0, 1$.

– Višeklasna logistička regresija:

$$\begin{aligned} h(\mathbf{x}; \mathbf{W}) &= f(\mathbf{s}) = \frac{1}{\mathbf{1}^\top \exp(\mathbf{s})} \exp(\mathbf{s}) = [P(y = k | \mathbf{x}, \mathbf{w})]_{k=1..C}^\top, \\ P(y | \mathbf{x}, \mathbf{w}) &= h_y(\mathbf{x}) = \prod_k h_k(\mathbf{x})^{\mathbb{I}[y=k]}, \\ L(y, h(\mathbf{x})) &= - \sum_k \mathbb{I}[y = k] \ln h(\mathbf{x})^{\mathbb{I}[y=k]}, \\ \nabla_{\mathbf{W}_{[k,:]}^\top} L(y, h_k(\mathbf{x})) &= (h_k(\mathbf{x}) - y_k) \phi(\mathbf{x}) \\ \nabla_{\mathbf{W}} L(y, h(\mathbf{x})) &= \phi(\mathbf{x})^\top (h(\mathbf{x}) - \mathbf{e}_y), \end{aligned}$$

gdje $y \in 1..C$, $h(\mathbf{x}) = [h_k(\mathbf{x})]_{k=1..C}^\top$, $h_i: \mathbb{R}^n \rightarrow \mathbb{R}$, a \mathbf{e}_k označava jednojedinичni vektor (vektor kanonske baze): $\mathbf{e}_k = [\mathbb{I}[i = k]]_{i=1..C}^\top$.

Funkcije gubitka su definirane kao negativni logaritmi izglednosti,

$L(y, h(\mathbf{x})) = -\ln P(y | \mathbf{x}, \mathbf{w})$, i konveksne su. Optimalne težine linearne regresije mogu se analitički izračunati, logistička regresija i višeklasna logistička regresija se obično uče optimizacijskim postupcima temeljenim na gradijentu.

Razdiobe $P(y | \mathbf{x}, \mathbf{w})$ poopćenih linearnih modela spadaju u **eksponencijalnu familiju razdioba**. Može se pokazati da je to jedina familija razdioba za koje

postoje konjugatne apriorne razdiobe, što pojednostavljuje računanje aposteriorne razdiobe (Murphy, 2012). Opći oblik ekponencijalne familije i više o njenima svojstvima i svojstvima poopćenih linearnih modela može se naći u Murphy (2012).

5. Duboko učenje i konvolucijske mreže

Klasični (plitki) modeli strojnog učenja (npr. poopćeni linearni modeli) oslanjaju se na kvalitetne značajke, tj. funkciju ϕ koja transformira ulazne primjere u vektore značajki. Za neke zadatke koji uključuju visokodimenzionalne primjere sa složenom strukturom (npr. slike, tekst i zvuk), ručno konstruiranje transformacije koja bi bila dovoljno dobra nije izvedivo, npr. jezgrene metode kod kojih se preslikavanje temelji na pretpostavci sličnosti primjera bliskih u ulaznom prostoru ne generaliziraju dobro. Kod **dubokog učenja** (LeCun et al., 2015; Goodfellow et al., 2016) transformacija ϕ se uči.

Odabirom

$$\phi(\mathbf{x}) = \phi(\mathbf{x}; \boldsymbol{\theta}_h) = f(\mathbf{W}_h \mathbf{x} + \mathbf{b}_h), \quad (5.1)$$

gdje je \mathbf{W}_h matrica težina, \mathbf{b}_h vektor pomaka, $\boldsymbol{\theta}_h = (\mathbf{W}_h, \mathbf{b}_h)$ a f nelinearna prijenosna (aktivacijska) funkcija koja se primjenjuje na svaki element vektora posebno, dobiva se jednostavna **umjetna neuronska mreža** (ovdje će se koristiti kraći nazivi: *neuronska mreža* ili *mreža*) s jednim **skrivenim slojem** kojemu odgovara funkcija ϕ . Ako to uvrstimo u jednadžbu poopćenog linearnog modela (jednadžba (4.13)):

$$h(\mathbf{x}; \boldsymbol{\theta}) = f(\mathbf{w}^T f(\mathbf{W}_h \mathbf{x} + \mathbf{b}_h) + b), \quad (5.2)$$

ili, ako je izlaz vektor,

$$h(\mathbf{x}; \boldsymbol{\theta}) = f(\mathbf{W}_o f(\mathbf{W}_h \mathbf{x} + \mathbf{b}_h) + \mathbf{b}_o), \quad (5.3)$$

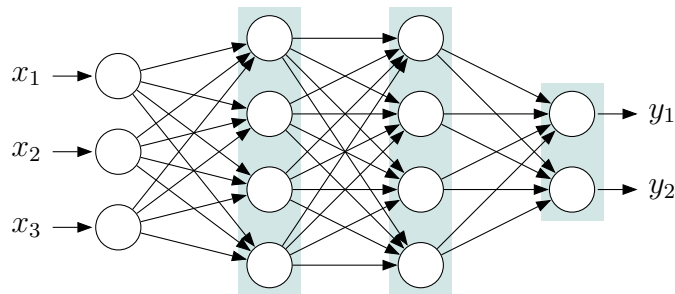
gdje $\boldsymbol{\theta} = (\mathbf{W}_h, \mathbf{b}_h, \mathbf{W}_o, \mathbf{b}_o)$. Jedinice neuronske mreže kojima odgovaraju operacije oblika $\mathbf{x} \mapsto f(\mathbf{w}_i^T \mathbf{x} + b_i)$ nazivaju se **umjetni neuroni**.

Za razliku od modela opisanih u odjeljku 4.7, za ovakav i dublje modele opisane u sljedećim odjeljcima ciljna funkcija nije konveksna (ni unimodalna) pa nije garantirano da će postupak učenja pronaći dobru hipotezu. Empirijski rezultati ipak pokazuju da duboke mreže uz neke prilagodbe uspješno uče i generaliziraju. Algoritmi koji se koriste za učenje modela dubokog učenja temelje se na gradijentnom spustu. Oni su opisani u odjeljku 5.3.

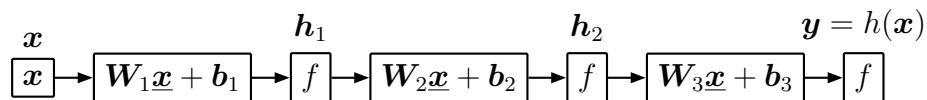
5.1. Duboke unaprijedne mreže

Može se pokazati da model mreže s jednim skrivenim slojem opisan jednačom (5.3), ako je dimenzija skrivenog sloja dovoljno velika, može s proizvoljno malom greškom aproksimirati svaku neprekinutu funkciju kojoj je domena konveksni podskup od \mathbb{R} . O tome govori teorem o univerzalnoj aproksimaciji (Cybenko, 1989; Leshno et al., 1993). Aktivacijska funkcija mora biti nelinearna jer kompozicija linearnih funkcija je linearna funkcija. Teorem o univerzalnoj aproksimaciji ne govori o tome hoće li takav model generalizirati. Dodavanjem jedinica u skriveni sloj povećava se kapacitet modela.

Obična neuronska mreža može imati više skrivenih slojeva, što se može prikazati kao na slici 5.1 ili slici 5.2. Povećavanjem broja skrivenih slojeva svaka jedinica u nekom sloju može koristiti izlaze svih jedinica prethodnog sloja kao značajke, što mreži omogućuje da, u odnosu na mrežu s 1 skrivenim slojem, s manje jedinica modelira funkcije u kojima postoje uzorci koji se ponavljaju i imaju hijerarhijsku strukturu (Goodfellow et al., 2016). Niži slojevi služe višim slojevima kao značajke transformiranjem kojih se dobivaju značajke više razine.



Slika 5.1: Prikaz primjera troslojne mreže. Svakom bridu odgovara jedna težina (pomaci nisu prikazani). Slojevi su označeni plavim četverokutima. Čvorovi koji su unutar slojeva mreže predstavljaju umjetne neurone. Slika je napravljena na temelju <http://www.texample.net/tikz/examples/neural-network/>.



Slika 5.2: Prikaz troslojne mreže kao računskog grafa. gdje čvorovi predstavljaju funkcije s parametrima, a bridovi podatke (vektore) čije su oznake prikazane uz neke od čvorova iz kojih izlaze. Funkcije su označene oznakom funkcije (aktivacijska funkcija) ili definicijom (afina transformacija). Ulaz sloja označen je s \underline{x} , a oznake varijabli koje pripadaju čvorovima (ulaz u ulaznom čvoru i parametri u čvorovima afine transformacije) nisu podvučene.

Kao prijenosna funkcija skrivenih slojeva često se koristi **zglobnica** (ReLU, engl. rectified linear unit) $\text{ReLU}(x) = \max(0, s)$, a prije su češće bile korištene **logistička funkcija**, $\sigma(s) = \frac{\exp(s)}{1+\exp(s)}$, i **tangens hiperbolni**, $\tanh(x) = \frac{\exp(s)-\exp(-s)}{\exp(s)+\exp(-s)}$. U izlaznom sloju obično se koriste funkcije korištene kod poopcenih linearnih modela (odjeljak 4.7) – identitet za regresiju, logistička funkcija za binarnu klasifikaciju, a **softmax**, $\text{softmax}(s) = \frac{1}{\mathbf{1}^\top \exp(s)} \exp(s)$, koji kao izlaz daje normalizirani vektor koji predstavlja razdiobu, za višeklasnu klasifikaciju.

Dosad opisivane mreže nazivaju se **unaprijedne mreže** (engl. *feedforward networks*) zato što se pri izračunu informacija propagira od ulaza prema izlazu, bez povratnih veza. Za mrežu kažemo da je duboka ako ima veći broj slojeva. Struktura duboke unaprijedne mreže ne mora se sastojati samo od niza afinih transformacija i nelinearnosti. Općenito, mrežu možemo predstaviti **računskim grafom**, tj. usmjerenim acikličkim grafom kod kojega čvorovi predstavljaju proizvoljne računske operacije i njihove izlaze, a bridovi označavaju izlazi kojih čvorova su ulazi kojih čvorova. Takav graf se naziva **računski graf**. Čvorovi koji nemaju roditelje su varijable koje čine ulazi i parametri. Parametri se mogu dijeliti, tj. mogu biti ulaz većem broju čvorova. Čvorovi koji nemaju djecu su izlazi (ili gubitak). Na slici 5.2 je prikazan jedan takav graf, ali na njemu, radi sažetosti, parametri nemaju zasebne čvorove, nego su označeni unutar čvorova koji o njima ovise.

5.2. Konvolucijske mreže

5.3. Učenje

5.3.1. Algoritam propagacije pogreške unatrag

5.3.2. Optimizacijski algoritmi

5.3.3. Algoritam propagacija pogreške unatrag

5.3.4. Isključivanje neurona - dropout

5.3.5. Normalizacija po grupama

6. Procjenjivanje nesigurnosti

6.1. Aleatorna i epistemička nesigurnost

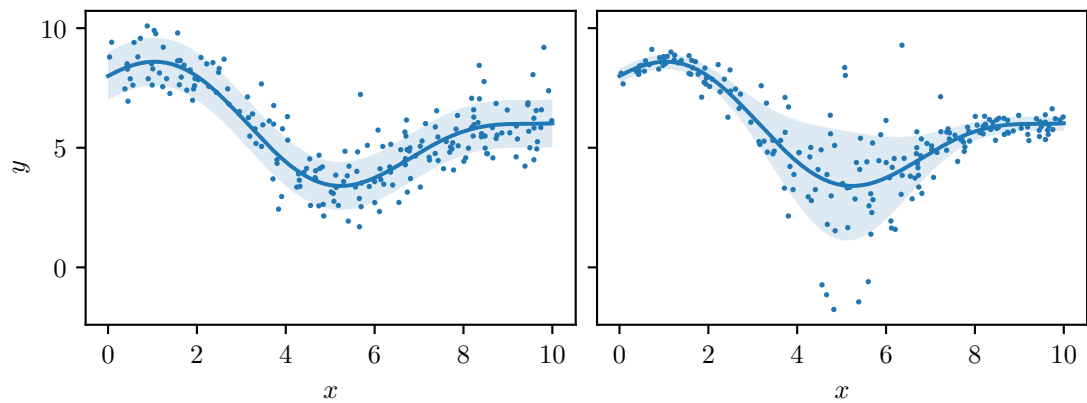
Kod bayesovskih modela nesigurnost zaključivanja izražava se razdiobom po vrijednostima varijable čija vrijednost se procjenjuje, a može se izraziti i entropijom ili varijancom kada je prikladno.

Postoje različiti izvori nesigurnosti (C. Kennedy i O'Hagan, 2002), ali nesigurnost općenito možemo podijeliti na dvije vrste: **aleatornu nesigurnost** i **epistemičku nesigurnost** (Kiureghian i Ditlevsen, 2009). 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 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.

6.2. Homoskedastička i heteroskedastička nesigurnost

Aleatorna nesigurnost može biti homoskedastička i heteroskedastička.



Slika 6.1: 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 $s(x) = 1$.

TODO: homoskedastička, heteroskedastička nesigurnost

TODO: model uncertainty, predictive uncertainty Gal-thesis 1.2

5.33331pt

11.74983pt

small 10.95pt

footnotesize 10.0pt

412.56497pt

7. Bayesovske neuronske mreže

8. Procenjivanje nesigurnosti kod konvolucijskih mreža

9. Eksperimentalni rezultati

9.1. Programska izvedba

9.2. Skupovi podataka

10. Zaključak

Zaključak.

LITERATURA

Ethem Alpaydin. **Introduction to Machine Learning**. 2014.

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>.

Anselm Blumer, Andrzej Ehrenfeucht, David Haussler, i Manfred K. Warmuth. Occam's razor. **Inf. Process. Lett.**, 24(6):377–380, Travanj 1987. ISSN 0020-0190. doi: 10.1016/0020-0190(87)90114-1. URL [http://dx.doi.org/10.1016/0020-0190\(87\)90114-1](http://dx.doi.org/10.1016/0020-0190(87)90114-1).

Anselm Blumer, A. Ehrenfeucht, David Haussler, i Manfred K. Warmuth. Learnability and the vapnik-chervonenkis dimension. **J. ACM**, 36(4):929–965, Listopad 1989. ISSN 0004-5411. doi: 10.1145/76359.76371. URL <http://doi.acm.org/10.1145/76359.76371>.

Marc C. Kennedy i Anthony O'Hagan. Bayesian calibration of computer models. 2002.

G. Cybenko. Approximation by superpositions of a sigmoidal function. **Mathematics of Control, Signals, and Systems (MCSS)**, stranice 303–314, 1989. ISSN 0932-4194. doi: 10.1007/BF02551274. URL <http://dx.doi.org/10.1007/BF02551274>.

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>.

Peter Grünwald. A tutorial introduction to the minimum description length principle. U **Advances in Minimum Description Length: Theory and Applications**, 2005.

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.

Yann LeCun, Yoshua Bengio, i Geoffrey E. Hinton. Deep learning. **Nature**, 521 (7553):436–444, 2015. doi: 10.1038/nature14539. URL <https://doi.org/10.1038/nature14539>.

Moshe Leshno, Vladimir Ya. Lin, Allan Pinkus, i Shimon Schocken. Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. **Neural Networks**, stranice 861–867, 1993. URL <http://dblp.uni-trier.de/db/journals/nn/nn6.html#LeshnoLPS93>.

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. Strojno učenje: 7. logistička regresija ii, 2017. URL http://www.fer.unizg.hr/_download/repository/SU-2017-07-LogistickaRegresija2.pdf.

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/>.

Samuel Rathmanner i Marcus Hutter. A philosophical treatise of universal induction. **CoRR**, abs/1105.5721, 2011. URL <http://arxiv.org/abs/1105.5721>.

Xitong Yang. Understanding the Variational Lower Bound, 2017. URL <http://legacydirs.umiacs.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.