

Tehtävä 1. Naivi-Bayes ja SPAM (2 pistettä).

Naivi-Bayes-mallia käytetään roskapostin suodattamiseen. Mallin todennäköisyysparametrien oppimista varten on kerätty valmiiksi luokiteltuja sähköpostiviestejä. Seuraavassa on joidenkin sanojen esiintymismäärät kummassakin luokassa (spam ja ham):

sana	spam	ham
Bayesin	0	5
kaavaa	2	8
me	9	10
rakastamme	11	1
<i>yht.</i>	9 000	12 000

Arvioi ehdolliset todennäköisyydet $P(\text{SANA}_i = s \mid \text{spam})$ ja $P(\text{SANA}_i = s \mid \text{ham})$, missä $s \in \{\text{kaavaa, Bayesin, rakastamme, me}\}$.

Laske näitä hyödyntäen seuraavat todennäköisyydet:

- i. (1 piste) $P(\text{spam} \mid \text{me})$
- ii. (1 piste) $P(\text{spam} \mid \text{me, rakastamme, Bayesin, kaavaa})$

Käytä prioritodennäköisyytenä arvoa $P(\text{spam}) = 0.5$.

(*Vihje:* Pienten todennäköisyyksien kohdalla on syytä käyttää jotakin alarajaa, esim. 0.0001. Muista että jos saat laskettua osamäärän $Odds = P(\text{spam} \mid \text{evidenssi})/P(\text{ham} \mid \text{evidenssi})$, saat todennäköisyyden kaavalla $P(\text{spam} \mid \text{evidenssi}) = Odds/(1 + Odds)$.)

Tehtävä 2. Numeroiden luokittelu neuroverkolla (2 pistettä).

Toteuta valmiiseen Java-ohjelmarunkoon, tai kokonaan alusta, perseptronialgoritmi numeroiden tunnistamiseen. Tiedosto `mnist-x.data` sisältää yhteensä 6000 kuvaa, yksi kuva jokaisella tiedoston rivillä, ja jokainen kuva on 28 x 28 pikseliä (eli jokaisella rivillä siis $28 \times 28 = 784$ arvoa). Jokainen pikseli on joko musta (-1) tai valkoinen (1). Tiedosto `mnist-y.data` sisältää näitä kuvia vastaavat luokat (0-9).

a) Suorita Java-pohja kerran ja varmista, että projektin runkoon ilmestyy tiedosto `test100.bmp`, jossa on sata ensimmäistä numeroa suuruusjärjestyksessä. Tällä tavalla verifioidaan että datan lukeminen on onnistunut.

b) (1 piste) Muokkaa perseptronialgoritmia (`Perseptroni`-luokan `train()`-metodi) siten, että se oppii erottamaan numeroa '3' (`targetChar`) esittävät kuvat numeroa '5' (`oppositeChar`) esittävistä kuvista. (Aseta luokkamuuttujan arvoksi 1, kun kuva esittää kolmosta, ja -1 jos se esittää vitosta.) Minkä luokitteluvirheen saat aikaiseksi?

c) (1 piste) Kokeile eri numeroparien erottelua (muitakin kuin 3 vs 5). Mitkä numerot on helpoin erottaa toisistaan, mitkä vaikein?

Tehtävä 3. Numeroiden luokittelu lähimmän naapurin luokittimella (1 piste).

Korvaa edellisen tehtävän ohjelmassa perseptroniluokittelija lähimmän naapurin luokittelijalla (tai koodaa kokonaan alusta).

Luokittelija etsii siis kullekin testiaineiston esimerkille, X , sitä vastaavan lähimmän opetusaineiston esimerkin, X^{train} , ja palauttaa sitä vastaavan luokan Y^{train} . Huomaa, että toisin kuin perseptroniluokittinta, lähimmän naapurin luokittinta ei varsinaisesti tarvitse opettaa vaan kaikki työ tehdään luokitteluvaiheessa.

Testaa taas luokittelijaasi samoilla opetus- ja testijoukoilla kuin edellisessä tehtävässä. Voit nyt ottaa mukaan kaikki luokat (numerot 0-9). Huomaa, että luokittelu 10:en luokkaan on vaikeampaa kuin luokittelu kahteen luokkaan, joten luokittelutarkkuus luultavasti alenee.