

Nulta laboratorijska vježba iz Dubokog učenja

Podsjetnik na logističku regresiju, gradijentni spust, Python i Numpy

U ovoj vježbi ćemo ponoviti osnove strojnog učenja, podsjetiti se na optimizaciju gradijentnim spustom te upoznati se s mogućnostima Pythona, Numpyja i Matplotliba. Razvit ćemo pythonske implementacije binarne i višerazredne logističke regresije utemeljene na Numpyjevim primitivima. Ispitati ćemo svojstva naše implementacije analizom numeričkih pokazatelja uspješnosti klasifikacije te iscrtavanjem funkcije odluke primjenom Matplotliba. Cilj vježbe je ponoviti osnovna znanja s prethodnih kolegija te razviti intuiciju koja će nam biti dragocjena pri oblikovanju i debuggiranju dubokih modela strojnog učenja.

U okviru vježbe razvit ćemo tri pythonska modula: `data`, `binlogreg` i `logreg`. Modul `data` će sadržavati operacije vezane uz stvaranje i iscrtavanje skupa slučajnih 2D podataka te operacije vezane uz evaluaciju klasifikacijske performanse. Taj modul ćemo koristiti i u prvoj vježbi gdje ćemo na istim podatcima usporediti duboko učenje sa SVM-om kao prvim izborom u mnogim klasifikacijskim problemima. Preostala dva modula će sadržavati funkcije i testove za binarnu i višerazrednu logističku regresiju.

0a. Uvodne napomene o parcijalnim derivacijama vektorskih funkcija

Deriviranje vektorskih i skalarnih funkcija više varijabli je ključni koncepti za shvaćanje gradijentne optimizacije parametara modela strojnog učenja i zato ga moramo dobro poznavati. Po definiciji, **parcijalna derivacija** funkcije više varijabli jest derivacija s obzirom na jednu od tih varijabli pod pretpostavkom da su sve ostale varijable konstantne. Parcijalne derivacije vektorskih funkcija možemo računati odvojeno za svaku komponentu funkcije.

Ilustrirajmo ove pojmove na primjeru sljedeće vektorske funkcije vektorske varijable:

$$\mathbf{f}(\mathbf{x}) = \begin{bmatrix} x_1^2 + x_2 x_3 \\ x_1 + x_2 + x_3 \end{bmatrix}$$

Parcijalna derivacija prve komponente funkcije \mathbf{f} po drugoj komponenti argumenta \mathbf{x} odgovara izrazu: $\partial f_1 / \partial x_2 = x_3$. Parcijalna derivacija druge komponente funkcije po prvoj komponenti argumenta odgovara izrazu: $\partial f_2 / \partial x_1 = 1$.

Jakobijeva matrica

Skup svih parcijalnih derivacija vektorske funkcije vektorske varijable možemo izraziti **Jakobijevom matricom** (ili, kraće, Jakobijanom). Retci Jakobijana odgovaraju komponentama vektorske funkcije, dok stupci odgovaraju komponentama nezavisne varijable. Tako npr. Jakobijan skalarne funkcije D-dimenzionalne vektorske varijable ima jedan redak i D stupaca, dok Jakobijan funkcije sinus ima jedan redak i jedan stupac. Po istoj logici, Jakobijan funkcije \mathbf{f} imao bi dva stupca i tri retka:

$$\mathbf{f}'(\mathbf{x}) = \frac{d\mathbf{f}}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \frac{\partial f_1}{\partial x_3} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \frac{\partial f_2}{\partial x_3} \end{bmatrix} = \begin{bmatrix} 2x_1 & x_3 & x_2 \\ 1 & 1 & 1 \end{bmatrix}$$

Često će nam biti interesantne matrice parcijalnih derivacija obzirom na samo neke od nezavisnih varijabli promatrane funkcije. Primjerice, ako kažemo da vektor \mathbf{q} sadrži prve dvije komponente vektora \mathbf{x} (tj. da vrijedi $\mathbf{q} = [x_1 x_2]^\top$), onda možemo pisati:

$$\mathbf{f}'(\mathbf{q}) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2x_1 & x_3 \\ 1 & 1 \end{bmatrix}$$

Kompozicije funkcija i pravilo ulančavanja

Prepostavimo sada da naša funkcija \mathbf{f} nema izravan kontakt s nezavisnom varijablom, nego da na nju djeluje posredno, preko funkcije \mathbf{g} . Rezultantna funkcija \mathbf{F} odgovara kompoziciji funkcija \mathbf{f} i \mathbf{g} :

$$\mathbf{F}(\mathbf{x}) = (\mathbf{f} \circ \mathbf{g})(\mathbf{x}) = \mathbf{f}(\mathbf{g}(\mathbf{x}))$$

Neka je funkcija \mathbf{g} zadana kako slijedi:

$$\mathbf{g}(\mathbf{x}) = \begin{bmatrix} \sin(x_1) \\ x_2^3 + x_1 x_2 \\ x_3 \end{bmatrix}$$

Prepostavimo da želimo odrediti derivaciju kompozicije funkcija \mathbf{F} . U pomoć nam priskače pravilo **ulančavanja** vektorskih funkcija koje u općenitom obliku glasi:

$$\mathbf{F}'(\mathbf{x}) = \mathbf{f}'(\mathbf{g}(\mathbf{x})) \cdot \mathbf{g}'(\mathbf{x})$$

Kad pravilo ulančavanja primijenimo na naš konkretni primjer, Jakobijevu matricu funkcije \mathbf{F} dobivamo sljedećim izvodom:

$$\begin{aligned} \mathbf{F}'(\mathbf{x}) &= \mathbf{f}'(\mathbf{g}(\mathbf{x})) \cdot \mathbf{g}'(\mathbf{x}) = \begin{bmatrix} 2g_1(\mathbf{x}) & g_3(\mathbf{x}) & g_2(\mathbf{x}) \\ 1 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} \cos(x_1) & 0 & 0 \\ x_2 & x_1 + 3x_2^2 & 0 \\ 0 & 0 & 1 \end{bmatrix} \\ \mathbf{F}'(\mathbf{x}) &= \begin{bmatrix} 2 \sin(x_1) & x_3 & x_2^3 + x_1 x_2 \\ 1 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} \cos(x_1) & 0 & 0 \\ x_2 & x_1 + 3x_2^2 & 0 \\ 0 & 0 & 1 \end{bmatrix} \\ \mathbf{F}'(\mathbf{x}) &= \begin{bmatrix} 2 \sin(x_1) \cos(x_1) + x_2 x_3 & x_1 x_3 + 3x_2^2 x_3 & x_2^3 + x_1 x_2 \\ \cos(x_1) + x_2 & x_1 + 3x_2^2 & 1 \end{bmatrix} \end{aligned}$$

Optimizacija skalarne funkcije više varijabli

U strojnem učenju posebno često analiziramo svojstva *skalarnih* funkcija više varijabli, a tipičan primjer je funkcija gubitka koju želimo optimirati. Ako je funkcija više varijabli skalarna (tj. ima samo jednu komponentu), tada *gradijent* te funkcije sadrži parcijalne derivacije po svim varijablama. Ako gradijent promatramo kao "obični" stupčani vektor (konvencije se razlikuju, ali to je najčešće slučaj), gradijent će odgovarati transponiranoj Jakobijevoj matrići:

$$\nabla f(\mathbf{x}) = \frac{df(\mathbf{x})}{d\mathbf{x}}^\top.$$

Ako funkciju f aproksimiramo Taylorovim razvojem prvog reda, dobivamo aproksimaciju:

$$f(\mathbf{x} + \Delta\mathbf{x}) = f(\mathbf{x}) + \frac{df(\mathbf{x})}{d\mathbf{x}} \Delta\mathbf{x} = f(\mathbf{x}) + \nabla f(\mathbf{x})^\top \Delta\mathbf{x}.$$

Pogledajmo što se zbiva kad u tu aproksimaciju uvrstimo pomak u smjeru negativnog gradijenta:

$$\Delta\mathbf{x} = -\epsilon \cdot \mathbf{g}, \quad \mathbf{g} = \nabla f(\mathbf{x}).$$

Ako Taylorova aproksimacija vrijedi, vrijednost funkcije f bi trebala opadati:

$$f(\mathbf{x} + \epsilon \cdot \mathbf{g}) = f(\mathbf{x}) - \epsilon \cdot \mathbf{g}^\top \mathbf{g}.$$

Uz malo sreće, iterativnim primjenjivanjem pomaka u smjeru negativnog gradijenta došli bismo do lokalnog minimuma funkcije f . To je najjednostavnija metoda gradijentne optimizacije, a poznata je pod nazivom gradijentni spust. U praksi, aproksimacija će biti tim točnija što je hiper parametar ϵ manji. Međutim, nije dobro da ϵ pude premali kako učenje ne bi previše trajalo. Zbog toga su smisljene brojne poboljšane metode, a neke od njih ćemo upoznati na predavanjima.

Ob. Uvodne napomene o višerazrednoj Logističkoj regresiji

Višerazredna logistička regresija [1] je probabilistički diskriminativni klasifikacijski model pripadnosti vektorskog podatka \mathbf{x} razredima $c_j, j = 0, 1, \dots, C - 1$. Radi se o generaliziranom linearnom modelu koji afinu transformaciju podataka spljošćuje (engl. squash) na vjerojatnosni interval $[0,1]$. Ishod klasifikacije formalno prikazujemo slučajnom varijablu Y za koju vrijedi $\sum_j P(Y = c_j | \mathbf{x}) = 1 \forall \mathbf{x}$.

U binarnom slučaju podatak može pripadati jednom od samo dvaju razreda ($C = 2$). Parametri modela tada su vektor \mathbf{w} te pomak b , a rezultat modela su aposteriorne vjerojatnosti razreda c_0 i c_1 :

$$P(Y = c_1 | \mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x} + b), \text{ gdje je } \sigma(s) = e^s / (1 + e^s)$$

$$P(Y = c_0 | \mathbf{x}) = 1 - P(c_1 | \mathbf{x})$$

Pri tome podatci \mathbf{x} imaju dimenzije $D \times 1$, parametar \mathbf{w} ima dimenzije $D \times 1$, dok je parametar b skalar.

Kad podatke klasificiramo u C razreda, model možemo sažeto formulirati sljedećim skupom jednadžbi (funkciju softmax čemo objasniti u nastavku):

$$P(Y = c_j | \mathbf{x}) = \text{softmax}(j, \mathbf{W} \cdot \mathbf{x} + \mathbf{b}), j = 0, 1, \dots, C - 1.$$

Matrice \mathbf{W} i \mathbf{b} sadrže parametre postupka, a njihove dimenzije su $C \times D$ odnosno $C \times 1$.

Pri tome D predstavlja dimenzionalnost podataka, dok je C broj razreda. Vektor $\mathbf{s} = \mathbf{W} \cdot \mathbf{x} + \mathbf{b}$ ima dimenzije $C \times 1$ te sadrži tzv. klasifikacijske mjere razreda (možemo koristiti i prijevod klasifikacijski rezultat, engl. classification scores). Klasifikacijske mjere pokazuju koliko je odgovarajući razred vjerojatniji (odnosno manje vjerojatan) u odnosu na druge razrede. Konačno, aposteriorne vjerojatnosti razreda dobivamo funkcijom softmax čija j -ta komponenta odgovara izrazu:

$$\text{softmax}(j, \mathbf{s}) = e^{s_j} / \sum_k e^{s_k}.$$

Ako detaljnije razmotrimo definiciju funkcije softmax, primjetit ćemo da su klasifikacijske mjere redundantne: kad svakoj od njih dodamo istu konstantu funkcija softmax na izlazu daje istu razdiobu. Iz toga slijedi da jednu od C klasifikacijskih mjer bez smanjenja općenitosti možemo fiksirati postavljanjem odgovarajućeg retka parametara na nulu ($\mathbf{W}_{j,:} = b_j = 0$). Tako normalizirane klasifikacijske mjeru mogu se interpretirati kao logaritmi omjera **šansi** promatranog razreda u odnosu na fiksirani razred. U praksi se pokazuje da postupak jednako dobro konvergira i ako optimiramo sve retke \mathbf{W} i \mathbf{b} . Stoga u implementaciji često optimiramo sve parametre jer se tako dobiva nešto jednostavniji programski kod. Međutim, redundantnost ulaza softmaxa u pravilu koristimo kako bismo izbjegli da eksponenciranje velikih klasifikacijskih mjer rezultira numeričkim preljevom. Primjer takve implementacije donosimo u nastavku:

```
# stabilni softmax
def stable_softmax(x):
    exp_x_shifted = np.exp(x - np.max(x))
    probs = exp_x_shifted / np.sum(exp_x_shifted)
    return probs
```

Parametre logističke regresije \mathbf{W} i \mathbf{b} određujemo optimiranjem prikladne funkcije gubitka $L(\mathbf{W}, \mathbf{b} | \mathbf{x}_i, y_i)$ na skupu za učenje. Skup za učenje tipično se sastoji od parova podataka i odgovarajućih oznaka razreda $\{(\mathbf{x}_i, y_i), i = 1, 2, \dots, N\}$. Ovo je najjednostavnije pojmiti na problemu klasifikacije 2D podataka u C različitim razreda. U tom slučaju \mathbf{x}_i su elementi ravnine, dok su y_i elementi cijelobrojnog intervala između 0 i $C-1$. Gubitak logističke regresije odgovara negativnoj log-izglednosti parametara:

$$L(\mathbf{W}, \mathbf{b} | \mathbf{x}_i, y_i) = \sum_i -\log P(Y = y_i | \mathbf{x}_i)$$

Ovako izražen gubitak favorizira parametre \mathbf{W} i \mathbf{b} koji točnim oznakama podataka za učenje pridružuju što veće vjerojatnosti. Može se pokazati da takav gubitak odgovara **unakrsnoj entropiji** (napravite to za vježbu!). U praksi nas ništa ne sprječava da umjesto negativne log-izglednosti parametara optimiramo neki **drugi gubitak**. Međutim, za slučaj klasifikacije, tj. kada su predikcije kategoričke, negativna log-izglednost vrlo dobro funkcionira u praksi te je najdosljedniji izbor kad učimo na čvrstim oznakama (za meke oznake treba nam unakrsna entropija).

Nažalost, optimiranje parametara logističke regresije ne možemo provesti u zatvorenoj formi, nego moramo posegnuti za nekim iterativnim postupkom. Na sreću, funkcija gubitka L je konveksna [2]. To znači da gradijentne metode poput gradijentnog spusta sigurno neće zapeti u nekom lokalnom optimumu (iako mogu podbaciti zbog drugih razloga).

Funkciju gubitka L možemo optimirati raznim gradijentnim postupcima. Takvi postupci u svakoj iteraciji učenja pomicu parametre u smjeru negativnog gradijenta gubitka.

Period tijekom kojeg algoritam dobije na uvid sve podatke za učenje nazivamo epohom. Ovisno o težini problema, za konvergenciju su potrebne desetine, stotine, ili tisuće epoha.

Glavni izazov kod gradijentnih pristupa optimizacije predstavlja određivanje gradijenta gubitka po svim parametrima postupka. S obzirom na to da gubitak i-tog podatka u logističkoj regresiji možemo promatrati kao kompoziciju logaritmiranog softmаксa i afine transformacije podataka, gradijente određujemo primjenom pravila [ulančavanja](#).

Literatura:

1. Jan Šnajder, Bojana Dalbelo Bašić: Strojno učenje. FER, Zagreb.
2. <http://qwone.com/~jason/writing/convexLR.pdf>

0c. Gradijenti binarne logističke regresije

Ovdje ćemo skicirati optimirani postupak računanja parcijalnih derivacija funkcije gubitka binarne logističke regresije. Za početak, dogovorimo notaciju. Označimo gubitak i-tog podatka s $L_i(\mathbf{w}, \mathbf{b}|\mathbf{x}_i, y_i) = -\log P(Y = y_i|\mathbf{x}_i)$. Definirajmo parcijalnu derivaciju $\partial L_i / \partial \mathbf{w}$ i predstavimo je retčanim vektorom u kojem se nalaze gradijenti L_i po elementima parametra \mathbf{w} . Slično, neka parcijalna derivacija $\partial L_i / \partial b$ bude retčani vektor gradijenata L_i po parametru b .

Kako bismo omogućili kompaktniji zapis aposteriorne vjerojatnost i-tog razreda, možemo uvesti *proširenu* sigmoidalnu funkciju koja uzima u obzir točan razred podatka čiju klasifikacijsku mjeru transformiramo u vjerojatnost:

$$\sigma_P(y, s) = \begin{cases} \sigma(s), & y = c_1 \\ 1 - \sigma(s), & y = c_0 \end{cases} = \begin{cases} \frac{e^s}{1+e^s}, & y = c_1 \\ \frac{1}{1+e^s}, & y = c_0 \end{cases}$$

Sada gubitak L_i odgovara kompoziciji logaritmirane proširene sigmoide i afine redukcije i-tog podatka (podsetimo se, s_i označava skalarnu klasifikacijsku mjeru podatka \mathbf{x}_i):

$$L_i(\mathbf{w}, \mathbf{b}|s_i) = -\log \sigma_P(y_i, s_i)$$

$$s_i = \mathbf{w}^\top \mathbf{x}_i + b$$

Gradijente parametara stoga možemo izraziti prema pravilu ulančavanja kao [umnožak](#) matrica parcijalnih derivacija (tzv. [Jacobijskih](#) matrica) funkcija L_i i s_i :

$$\partial L_i / \partial \mathbf{w} = \partial L_i / \partial s_i \cdot \partial s_i / \partial \mathbf{w}$$

$$\partial L_i / \partial b = \partial L_i / \partial s_i \cdot \partial s_i / \partial b$$

Parcijalna derivacija funkcije gubitka L po linearном klasifikacijskom rezultatu s u podatku \mathbf{x}_i jest skalar $\partial L_i / \partial s_i$ kojeg promatramo kao Jakobijan dimenzija 1×1 .

Parcijalne derivacije klasifikacijske mjere podatka \mathbf{x}_i po elementima parametara \mathbf{w} odnosno b , su matrica $\partial s_i / \partial \mathbf{w}$ dimenzija $1 \times D$ odnosno skalar $\partial s_i / \partial b$ kojeg promatramo kao matricu 1×1 . Izrazimo prvo derivaciju gubitka po klasifikacijskoj mjeri tako da ovisnost o točnom razredu i-tog podatka izrazimo [Iversonovim](#) uglatim zagradama. Ako se dogovorimo da vrijednost izraza $\llbracket q \rrbracket$ iznosi 1 ako je q istinit a nula inače, traženu parcijalnu derivaciju možemo predstaviti sljedećim izrazom (provjerite to za vježbu!):

$$\partial L_i / \partial s_i = P(c_1 | \mathbf{x}_i) - \llbracket y_i = c_1 \rrbracket.$$

Parcijalna derivacija klasifikacijske mjere po parametrima \mathbf{w} i b za podatak \mathbf{x}_i sada je:

$$\frac{\partial s_i}{\partial \mathbf{w}} = \mathbf{x}_i^\top$$

$$\frac{\partial s_i}{\partial b} = 1.$$

Kad uvedemo pokratu $\mathbf{g}_s = [\partial L_i / \partial s_i]_{i=1}^N$ te uzmemu u obzir da ukupni gubitak zbraja doprinose po svim podatcima, dobivamo konačne izraze:

$$\frac{\partial L}{\partial \mathbf{w}} = 1/N \cdot \sum_i \frac{\partial L_i}{\partial \mathbf{w}} = 1/N \cdot \sum_i g_{si} \cdot \mathbf{x}_i^\top,$$

$$\frac{\partial L}{\partial b} = 1/N \cdot \sum_i \frac{\partial L_i}{\partial b} = 1/N \cdot \sum_i g_{si}.$$

Najbolja računska performansa danas se postiže izražavanjem petlji optimiranim operacijama nad matricama i vektorima. Zbog toga je jasno da ćemo u implementaciji računanje gradijenta pomaka izraziti pozivom odgovarajuće funkcije numeričke biblioteke. Manje je međutim jasno kako izbjegći pisanje petlje za računanje gradijenata težina, pa ćemo tom važnom detalju pokloniti više pažnje.

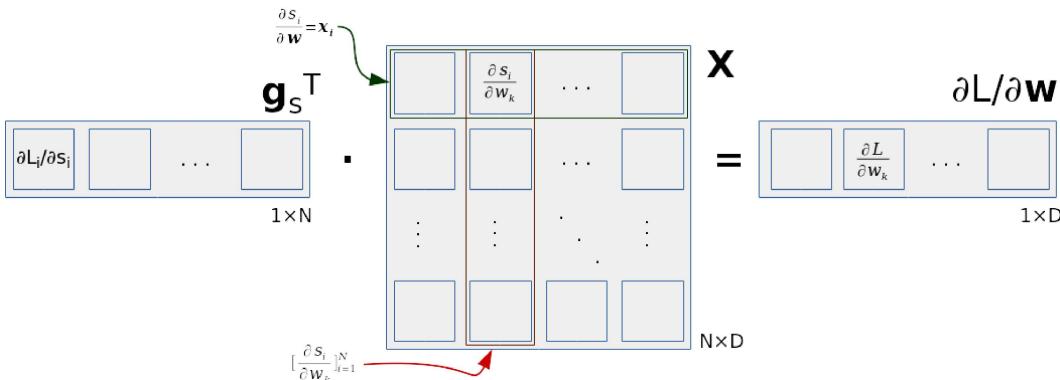
Naime, pažljivim promatranjem možemo uočiti da računanje parcijalne derivacije $\frac{\partial L}{\partial w}$ možemo vektorizirati tako da zbroj po podatcima predstavimo matričnim umnoškom. Vektore podataka ćemo organizirati u matricu \mathbf{X} dimenzija NxD. Retci matrice \mathbf{X} odgovaraju podatcima \mathbf{x}_i . Parcijalne derivacije gubitaka po klasifikacijskoj mjeri smjestit ćemo u stupčani vektor \mathbf{g}_s dimenzija Nx1. Kao što smo naveli ranije, vrijedi: $\mathbf{g}_s = [\partial L_i / \partial s_i]_{i=1}^N$. Lako se vidi da gradijent gubitka po k-toj težini odgovara skalarnom produktu vektora \mathbf{g}_s i k-tog stupca matrice \mathbf{X} :

$$\frac{\partial L_i}{\partial w_k} = 1/N \cdot \sum_i g_{si} \cdot x_{ik} = 1/N \cdot \mathbf{g}_s^\top \mathbf{X}_{:,k}.$$

Sada je jasno da parcijalne derivacije ukupnog gubitka po svim težinama možemo dobiti produktom vektora \mathbf{g}_s^\top i podatkovne matrice \mathbf{X} :

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{g}_s^\top \cdot \mathbf{X}.$$

Taj umnožak prikazan je na sljedećoj slici. Na slici se vidi da retci podatkovne matrice sadrže pojedinačne podatke te odgovaraju parcijalnoj derivaciji klasifikacijske mjere po vektoru težina za taj podatak. Stupci podatkovne matrice odgovaraju parcijalnim derivacijama po pojedinačnim komponentama vektora težina. Parcijalnu derivaciju ukupnog gubitka po k-toj težini dobivamo skalarnim produktom vektora \mathbf{g}_s i k-tog stupca matrice \mathbf{X} .



Napominjemo da postoji dobar razlog da matrica podataka ima strukturu NxD. U tom su slučaju komponente pojedinog podatka pohranjene na uzastopnim memorijskim lokacijama pa miješanje podataka i operacije nad grupama (engl. batch) uzrokuju manje promašaja priručne memorije. U programskoj implementaciji htjet ćemo dobiti gradiente grad_w koji po dimenzijama odgovaraju originalnim parametrima pa ćemo gornji izraz naprsto transponirati. U tom slučaju, poboljšane vrijednosti parametara dobivamo oduzimanjem gradijenata grad_w moduliranih hiperparametrom param_delta .

Od. Gradijenti višerazredne logističke regresije

Jednadžbe gradijenata parametara višerazredne logističke regresije vrlo su slične odgovarajućim gradijentima binarne logističke regresije. To nas ne treba čuditi jer se lako uviđa da je softmax popočenje sigmoidalne funkcije. Međutim, izražavanje tih gradijenata otežava struktura parametra \mathbf{W} . Da izbjegnemo poteškoće, odvojeno ćemo promatrati gradijente svakog od C redaka matrice \mathbf{W} , te usporedno s njima gradijente

svakog od C elemenata vektora \mathbf{b} . Označimo te retke odnosno elemente s $\mathbf{W}_{j:}$ i b_j , a njihove parcijalne derivacije s $\partial L / \partial \mathbf{W}_{j:}$ (matrica dimenzija 1xD) i $\partial L / \partial b_j$ (matrica dimenzija 1x1). Kao i ranije, gradijente je najlakše prikazati tako da ovisnost o točnom razredu i-tog podatka prikažemo Iversonovim uglatim zagradama. Prvo ćemo izraziti gradijent funkcije gubitka za i-ti podatak L_i po j-toj klasifikacijskoj mjeri (izvedite ovaj izraz za vježbu!).

$$\partial L_i / \partial s_{ij} = P(Y = c_j | \mathbf{x}_i) - [\![y_i = c_j]\!].$$

U izvedbi je ove parcijalne derivacije najpraktičnije sakupiti u matricu

$\mathbf{G}_s = [[\partial L_i / \partial s_{ij}]_{j=1}^C]_{i=1}^N$ dimenzija NxN. Kako bismo tu matricu efikasno izračunali, uvodimo matricu aposteriornih vjerojatnosti podataka $\mathbf{P}_{ij} = P(c_j | \mathbf{x}_i)$ te matricu vektorski kodiranih oznaka \mathbf{Y} :

$$\mathbf{Y}_{ij} = \begin{cases} 1, & y_i = c_j \\ 0, & y_i \neq j \end{cases}.$$

Sada matricu parcijalnih derivacija gubitka po klasifikacijskim mjerama dobivamo jednostavnim matričnim izrazom $\mathbf{G}_s = \mathbf{P} - \mathbf{Y}$.

Parcijalne derivacije klasifikacijskih mjera po parametrima iste su kao i kod binarne logističke regresije:

$$\partial s_{ij} / \partial \mathbf{W}_{j:} = \mathbf{x}_i^\top$$

$$\partial s_{ij} / \partial b_j = 1.$$

Kad ulančimo gradijente komponenata gubitka te uzmemu u obzir da ukupni gubitak zbraja doprinose po svim podatcima, dobivamo izraz kojeg možemo koristiti u optimizacijskom postupku:

$$\partial L / \partial \mathbf{W}_{j:} = 1/N \cdot \sum_i \partial L_i / \partial \mathbf{W}_{j:} = 1/N \cdot \sum_i \partial L_i / \partial s_{ij} \cdot \partial s_{ij} / \partial \mathbf{W}_{j:}$$

$$\partial L / \partial \mathbf{W}_{j:} = 1/N \cdot \sum_i G_{sij} \cdot \mathbf{x}_i^\top$$

$$\partial L / \partial b_j = 1/N \cdot \sum_i \partial L_i / \partial b_j = 1/N \cdot \sum_i G_{sij}.$$

Pažljivim promatranjem ovog izraza uočit ćemo ogromnu sličnost s odgovarajućim izrazom gradijenata parametara binarne logističke regresije. To nas ne treba čuditi, jer retci matrice težina \mathbf{W} u višerazrednom gubitku sudjeluju na potpuno isti način kao i težine \mathbf{w} u gubitku binarne logističke regresije. Komponente gradijenta $\partial L / \partial \mathbf{W}_{j:}$ odgovaraju skalarnom produktu j-tog stupca matrice \mathbf{G}_s sa stupcima podatkovne matrice \mathbf{X} . Stoga bismo, kao i kod binarne logističke regresije, sve gradijente j-tog retka matrice težina mogli izračunati umnoškom vektora $\mathbf{G}_{s:j}$ i matrice \mathbf{X} :

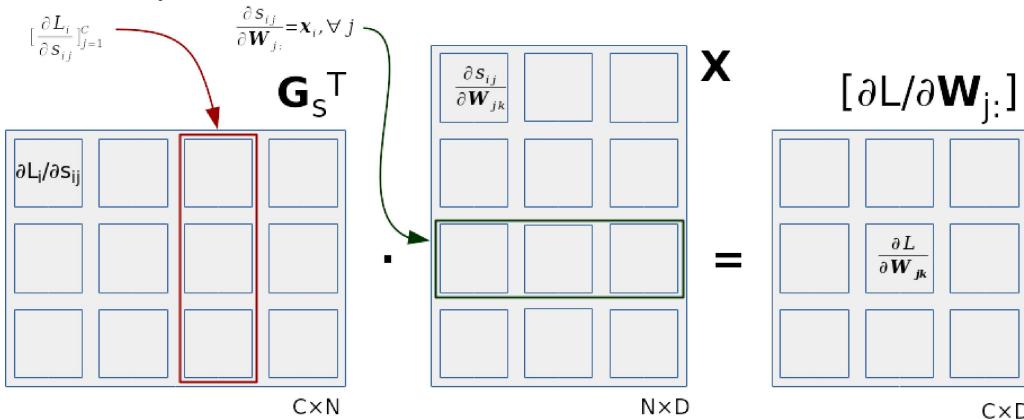
$$\partial L / \partial \mathbf{W}_{j:} = \mathbf{G}_{s:j} \cdot \mathbf{X}.$$

Međutim, pokazuje se da u praksi možemo i bolje od toga. Važno svojstvo našeg problema jest da su parcijalne derivacije klasifikacijskih mjera po odgovarajućim redcima matrice težina međusobno jednake: $\partial s_{ij} / \partial \mathbf{W}_{j:} = \mathbf{x}_i^\top \forall j$. Odatle slijedi da gradijente svih redaka matrice težina vrlo praktično možemo izračunati samo jednim matričnim umnoškom:

$$[\partial L / \partial \mathbf{W}_{j:}]_{j=1}^C = \mathbf{G}_s^\top \cdot \mathbf{X}$$

Prikazani izraz relativno je lako zapamtiti, ali iznimno teško objasniti reverse engineeringom. Često se može čuti objašnjenje kako do tog izraza dolazimo izravnom primjenom pravila ulančavanja. Međutim, takvo objašnjenje je netočno jer pravilom ulančavanja ne možemo računati parcijalne derivacije po matričnoj varijabli. Ako ustraјemo na takvom pristupu, morat ćemo ili izravnati matricu u vektor ili prihvati da Jakobijan u tom slučaju postaje tenzor trećeg reda. U oba slučaja složenost našeg algoritma bi se povećala jer ne bismo mogli iskoristiti inherentnu algebarsku strukturu problema, koju možemo jednostavno izraziti ovako: $\partial s_{ij} / \partial \mathbf{W}_{k:} = 0, j \neq k$. Stoga, izraz $\mathbf{G}_s^\top \cdot \mathbf{X}$ treba promatrati kao *komplikiranu optimizaciju*, a ne kao početnu točku za razumijevanje načina za određivanje gradijenata matrice težina višerazredne logističke regresije. Molimo studente koji prakticiraju nelinearno učenje da ovu netrivijalnu ali vrlo važnu činjenicu prihvate i zapamte što je moguće prije.

Optimirani izraz za računanje gradijenata višerazredne logističke regresije prikazan je na sljedećoj slici. Na slici vidimo da doprinos svakog podatka gradijentima težina odgovara vanjskom umnošku stupca matrice \mathbf{G}_s^T (odnosno, retka matrica \mathbf{G}_s) i odgovarajućeg retka podatkovne matrice. Kad bismo imali jedan podatak (indeks i) i kad bismo računali gradijente samo jednog retka matrice težina (indeks j), matrica \mathbf{G}_s svela bi se na element \mathbf{G}_{sij} , a matrica \mathbf{X} na redak \mathbf{x}_i .



U praksi, gradijente težina $[\frac{\partial L}{\partial \mathbf{W}_j}]_{j=1}^C$ računat ćemo jednim pozivom funkcije matrične biblioteke, npr. `grad_W = np.dot(GsT, x)`. Na taj način jednim pozivom rješavamo sljedeća dva zadatka:

- istovremeno računanje gradijenata svih redaka matrice \mathbf{W}
- akumuliranje doprinosa svih podataka \mathbf{x}_i ;

Vrlo sličan postupak imat ćemo i za gradijent pomaka \mathbf{b} , samo što ćemo umjesto matričnog množenja pozivati funkciju za istovremeno zbrajanje svih stupaca matrice, npr. `np.sum(GsT, axis=0)`.

Ovakve manevre vrlo rado ćemo koristiti u praksi kako programske petlje ne bi bile u programskom jeziku visoke razine nego u optimiranom bibliotečnom kodu koji može biti i do 100 puta brži od naivnog koda u C-u. Naime, ako je numpy preveden s [OpenBLAS-om](#), izvedba matričnog množenja će koristiti optimirani kod koji pazi na promašaje cachea i temelji se na ručno optimiranom [strojnom kodu](#) koji je prilagođen za konkretnu mikroarhitekturu. Nadalje, ako je OpenBLAS konfiguriran da koristi OpenMP, u naš algoritam biti će upregnute sve procesorske jezgre - a to bi rezultiralo dodatnim ubrzanjem. Neke konkretne brojeve možete pogledati [ovdje](#) (u tim eksperimentima OpenMP je bio isključen).

0e. Uvodne napomene vezane uz Python

Python je moderni dinamički jezik opće namjene koji je zbog svoje univerzalne primjenljivosti postao iznimno popularan u područjima vezanim uz umjetnu inteligenciju. Ako niste položili Skriptne jezike, predlažemo da samostalno proučite [skriptu](#) ili [knjigu](#). Koristit ćemo Python 3. Preporučamo vam da isprobate naprednu interaktivnu ljsku `ipython`.

Izvorni kod u Pythonu smještamo u datoteke koje nazivamo modulima. Svaki modul u sebi treba sadržavati kratki ispitni program koji testira funkcionalnost modula i ilustrira njegovo pravilno korištenje. Ispitne programe smještamo na kraj modula u tijelu uvjetne naredbe koja testira je li modul pokrenut kao glavni program. Na taj način omogućavamo provođenje testiranja jednostavnim izvršavanjem modula.

```
if __name__ == "__main__":
    # test, test, test!
```

Postoji više načina za debuggiranje u Pythonu. Jedan od najzgodnijih oslanja se na ugrađeni Pythonov debugger `pdb` kojeg možemo pozvati izravno iz koda. Da bismo to proveli, potrebno je prvo uključiti odgovarajući modul:

```
import pdb
```

Zatim, na mjestu gdje želimo prekinuti izvođenje (breakpoint!) navedemo poziv:

```
pdb.set_trace()
```

Nakon tog poziva Python će otvoriti interaktivnu **Ijusku** u kojoj će biti dostupna sva imena programa koja su bila vidljiva u trenutku poziva funkcije `set_trace`. Iz Ijuske možemo ispisivati objekte, pozivati proizvoljne funkcije (!) i izvršavati proizvoljne naredbe Pythona (!). Pored toga, na raspolaganju nam je i čitav niz specijalnih naredbi za debuggiranje. Popis tih naredbi dobivamo naredbom `help` (neke od važnijih naredbi su `up`, `down` i `continue`). Sličan učinak možemo postići i pozivom funkcije `IPython.embed()` (međutim, tu naredbe za debuggiranje nisu dostupne).

Na kraju ćemo pokazati kako debuggirati kod koji baca neočekivane iznimke na sljedećem minimalnom primjeru:

```
def proba():
    a=5
    raise ValueError("Iznimka!")

if __name__=="__main__":
    proba()
```

Testiranjem ovog programa vrlo lako se možemo uvjeriti da njegovo izvršavanje završava iznimkom `ValueError`. Malo teže je međutim doći do podrobnejih informacija o stanju programa na mjestu gdje se iznimka dogodila, poput npr. ispisivanja vrijednosti imena `a`. Problem možemo rješiti pokretanjem debuggera `pdb proba.py` te zadavanjem naredbe `continue` (skraćeno "c"). Debugger će početi izvoditi program te kad dođe do neuhaćene iznimke omogućiti će nam uvid u stanje programa preko tekstovnog korisničkog sučelja. Primjerice, naredba `print a` će ispisati vrijednost imena `a` neposredno prije bacanja iznimke. Do debuggera možemo doći i iz alternativne Ijuske `ipython` na način da nakon interaktivnog poziva čije izvršavanje završava neuhaćenom iznimkom zadamo naredbu `debug`.

0f. Uvodne napomene vezane uz Numpy

Numpy je popularna numerička biblioteka za Python. U primjerima ćemo koristiti standardnu pokratu `np` do koje dolazimo naredbom:

```
import numpy as np
```

Matplotlib je biblioteka koja se često kombinira s numpyjem, a služi za izradu raznih vrsta grafova u Pythonu. U primjerima ćemo koristiti standardnu pokratu `plt` do koje dolazimo naredbom:

```
import matplotlib.pyplot as plt
```

Važna prednosti numpyja jest mogućnost svođenja iterativnih postupaka na vektorske i matrične izraze. Za primjer, pogledajmo kako bismo u numpyju izrazili afinu transformaciju podataka koja predstavlja važan korak pri učenju logističke regresije. Neka su podatci zadani numpyjevom matricom `x` dimenzija $N \times D$ (svaki redak odgovara jednom D -dimenzijskom podatku), neka je transponirana matrica težina \mathbf{W}^T zadana numpyjevom matricom `w_t` dimenzija $D \times C$ te neka je transponirani vektor pomaka \mathbf{b}^T zadan numpyjevim vektorom `b_t` dimenzija $1 \times C$. Tada afinu transformaciju *svih* podataka dobivamo jednostavnim pozivom numpyjevih funkcija za matrično množenje i zbrajanje:

```
H = np.dot(x, w_t) + b_t
```

Obratimo međutim pozornost na jedan detalj koji bi nakon površnog pogleda mogao i promaknuti. Pažljivom analizom prikazane naredbe Pythona uočavamo da dimenzije pribrojnika ne odgovaraju: matrični umnožak `np.dot(x, w_t)` ima dimenzije $N \times C$ dok

vektor \mathbf{b}_t ima dimenziju $1 \times C$ ($N \neq 1$). Biblioteka numpy takva nesuglasja rješava pod pretpostavkom da vektor \mathbf{b}_t treba dodati *svakom* retku matričnog umnoška. Stoga će nakon izvršavanja naredbe dodjeljivanja ime \mathbf{H} uistinu referencirati matricu čiji retci odgovaraju afino transformiranim retcima podatkovne matrice. Kad se podsjetimo da se i-ti podatak nalazi u i-tom retku matrice podataka ($\mathbf{x}_i = \mathbf{X}_{i:}^T$), konačnu vrijednost i-tog retka matrice \mathbf{H} moći ćemo zapisati i sljedećim izrazom: $\mathbf{H}_{i:}^T = \mathbf{W} \cdot \mathbf{x}_i + \mathbf{b}$.

Upravo smo vidjeli da numpy po potrebi automatski *umnožava* operand manje dimenzionalnosti (engleski termin je broadcasting) kako bi došao do valjanog izraza. Ovakvo izražavanje je ispočetka vrlo teško prihvatiti, ali nakon nekog vremena postaje jasno da prednosti (sažetost i efikasnost) uvjerljivo nadjačavaju nedostatke.

1. Stvaranje umjetnog skupa 2D podataka

Rješenje ove vježbe slobodno preuzmite [ovdje](#).

U ovoj vježbi ćemo izvesti razred `Random2DGaussian` koji će nam omogućiti uzorkovanje 2D podataka generiranih slučajnom Gaussovom distribucijom. Konstruktor razreda treba stvarati parametre slučajne bivariatne Gaussove razdiobe (vektor μ i matricu Σ). Metoda `get_sample(n)` treba vratiti n slučajno uzorkovanih 2D podataka u numpyjevom polju dimenzija $N \times 2$. Položaj i varijanca razdiobe trebaju biti ograničeni fiksnim parametrima.

Pomoć:

- Definirajte raspon prihvatljivih vrijednosti za obje koordinate sredine razdiobe μ (npr. `minx=0, maxx=10, miny=0, maxy=10`); neka to budu podatkovni članovi razreda.
- Slučajno izaberite sredinu razdiobe μ prema uniformnoj razdiobi (`np.random.random_sample`).
- Slučajno odaberite svojstvene vrijednosti kovarijacijske matrice Σ prema uniformnoj razdiobi i organizirajte ih u dijagonalnu matricu D . Neka gornja granica svojstvenih vrijednosti ovisi o rasponu prihvatljivih vrijednosti parametra μ , npr.

```
eigvalx = (np.random.random_sample()*(maxx - minx)/5)**2
```

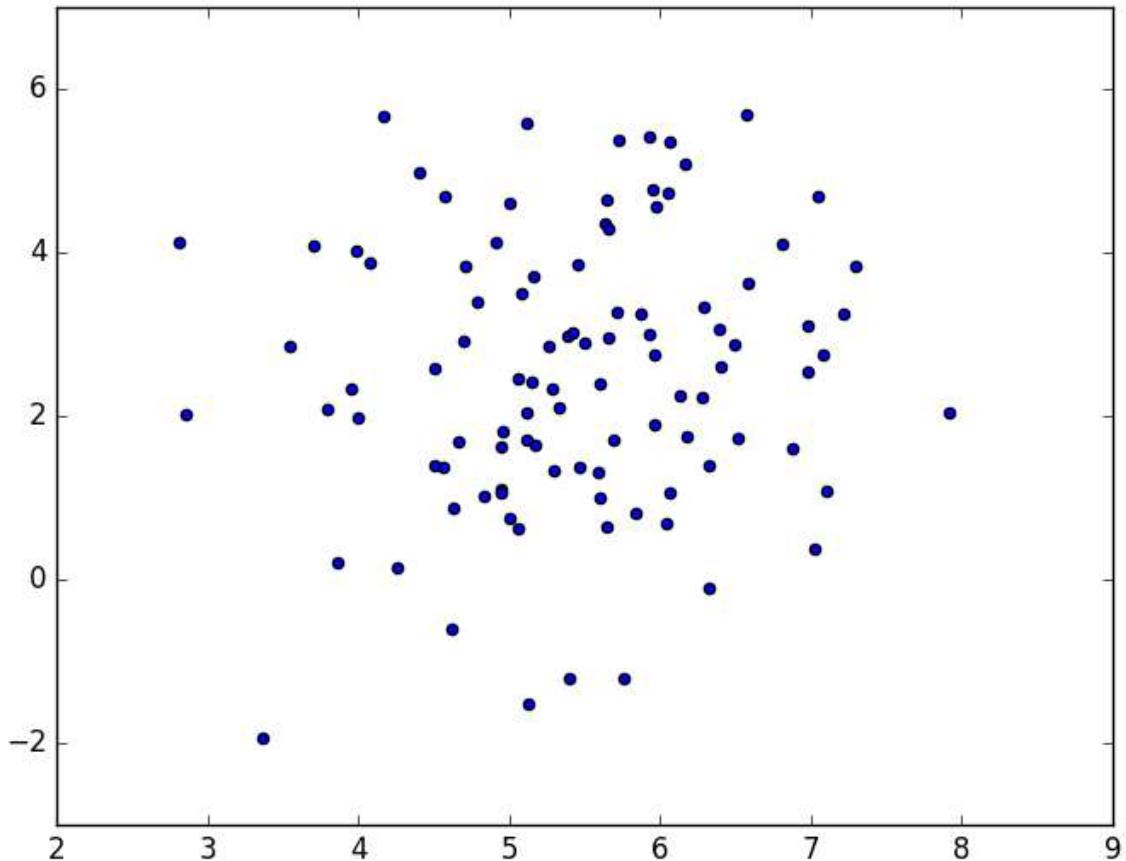
- Slučajno odaberite kut rotacije kovarijacijske matrice i na temelju njega formirajte rotacijsku matricu $R = \begin{bmatrix} \cos(\varphi) & -\sin(\varphi) \\ \sin(\varphi) & \cos(\varphi) \end{bmatrix}$.
- Odredite matricu Σ kao umnožak transponirane matrice R , matrice D i matrice R .
- Za uzorkovanje koristite funkciju `np.random.multivariate_normal`.
- Kako biste osigurali ponovljivost eksperimenata sa slučajno generiranim podatcima, numpyjev generator slučajnih brojeva slobodno inicijalizirajte na konstantu:

```
np.random.seed(100)
```

Ispravnost koda provjerite sljedećim ispitnim programom:

```
if __name__=="__main__":
    G=Random2DGaussian()
    X=G.get_sample(100)
    plt.scatter(X[:,0], X[:,1])
    plt.show()
```

Ovisno o stanju generatora slučajnih brojeva, vaš rezultat mogao bi izgledati kao na sljedećoj slici:



Ako je rezultat prihvatljiv, pohranite kod u datoteku data.py.

2. Učenje binarne logističke regresije gradijentnim spustom

Napišite funkciju `binlogreg_train` koja za zadani skup podataka za učenje optimizira parametre logističke regresije w i b . U vašem izvornom kodu koristite sljedeća imena:

- x : matrica podataka dimenzija $N \times D$;
- $y_$: vektor *točnih* razreda podataka dimenzija $N \times 1$ (koristimo ga tijekom učenja);
- y : vektor *predviđenih* razreda podataka dimenzija $N \times 1$ (koristimo ga tijekom ispitivanja performanse).

Neka funkcija ima sljedeće sučelje:

```
def binlogreg_train(X,Y_):
    """
    Argumenti
        X: podatci, np.array NxD
        Y_: indeksi razreda, np.array Nx1

    Povratne vrijednosti
        w, b: parametri logističke regresije
    """

```

Upute:

- Inicijalizirajte w prema normalnoj distribuciji $N(0,1)$ (`np.random.randn`) te b na nulu.
- Otvorite petlju gradijentnog spusta. Broj iteracija postupka zadajte hiperparametrom `param_niter`.

- Izračunajte klasifikacijske mjere svih podataka na način da petlje prepustite numpyju (`scores = np.dot(X, w) + b`).
- Izračunajte aposteriorne vjerojatnosti razreda u svim podatcima $P(c_1|\mathbf{X})$ na način da petlje prepustite numpyju
- Izračunajte gubitak $L(\mathbf{w}, \mathbf{b})$ na način da petlje prepustite numpyju.
- Izračunajte gradijente `grad_w` i `grad_b` na način da petlje prepustite numpyju.
- Izmjenite parametre u smjeru negativnog gradijenta (neka faktor pomaka `param_delta` bude hiperparametar postupka).

Struktura vaše izvedbe gradijentnog spusta trebala bi izgledati ovako:

```
# gradijentni spust (param_niter iteracija)
for i in range(param_niter):
    # klasifikacijske mjere
    scores = ...      # N x 1

    # vjerojatnosti razreda c_1
    probs = ...        # N x 1

    # gubitak
    loss = ...         # scalar

    # dijagnostički ispis
    if i % 10 == 0:
        print("iteration {}: loss {}".format(i, loss))

    # derivacije gubitka po klasifikacijskim mjerama
    dL_scores = ...      # N x 1

    # gradijenti parametara
    grad_w = ...        # D x 1
    grad_b = ...        # 1 x 1

    # poboljšani parametri
    w += -param_delta * grad_w
    b += -param_delta * grad_b
```

Napišite funkciju `binlogreg_classify` koja klasificira zadani skup podataka u skladu sa zadanim parametrima logističke regresije. Neka funkcija ima sljedeće sučelje:

```
...
Argumenti
    X:      podatci, np.array NxD
    w, b:  parametri logističke regresije

Povratne vrijednosti
    probs: vjerojatnosti razreda c1
...
```

Napišite funkciju `sample_gauss_2d(c, N)` koja stvara c slučajnih bivarijatnih Gaussovih razdioba, te iz svake od njih uzorkuje N podataka. Funkcija treba vratiti matricu x dimenzija (N·C)x2 čiji retci odgovaraju uzorkovanim podatcima te matricu točnih razreda y dimenzija (N·C)x2 čiji jedini stupac sadrži indeks razdiobe iz koje je uzorkovan odgovarajući podatak. Ako je i-ti redak matrice x uzorkovan iz razdiobe j, onda mora biti $Y[i, 0] = j$.

Napišite funkciju `eval_perf_binary(Y, Y_r)` koja na temelju predviđenih i točnih indeksa razreda određuje pokazatelje performanse binarne klasifikacije: točnost (engl. accuracy), preciznost (engl. precision) te odziv (engl. recall). Implementaciju te funkcije temeljite na brojnostima istinitih pozitiva (TP), lažnih pozitiva (FP), istinitih negativa (TN) i lažnih negativa (FN).

Napišite funkciju `eval_AP` koja računa prosječnu preciznost binarne klasifikacije. Neka funkcija na ulazu prima rangiranu listu točnih razreda \mathbf{Y}_r , koju dobivamo kad matricu točnih razreda \mathbf{Y} sortiramo prema aposteriornim vjerojatnostima odgovarajućih podataka $P(c_1|\mathbf{x})$. Rangiranu listu točnih razreda možete dobiti pozivom metode `argsort` numpyjevog polja. Prosječnu preciznost možete izračunati primjenom sljedećeg

izraza: $AP = \frac{\sum_{i=0}^{N-1} \text{Preciznost}(i) \cdot \mathbf{Y}_{r_i}}{\sum_{i=0}^{N-1} \mathbf{Y}_{r_i}}$. Pri tome $\text{Preciznost}(i)$ odgovara preciznosti u slučaju kad podatke s indeksom većim ili jednakim i pridružimo razredu c_1 , a podatke s indeksom manjim od i - razredu c_0 . Evo kako bi se trebala ponašati vaša funkcija:

```
>>> import data
>>> data.eval_AP([0,0,0,1,1,1])
1.0
>>> data.eval_AP([0,0,1,0,1,1])
0.9166666666666666
>>> data.eval_AP([0,1,0,1,0,1])
0.7555555555555555
>>> data.eval_AP([1,0,1,0,1,0])
0.5
```

Napišite ispitni kod za modul `binlogreg.py`. Formirajte skup za učenje pozivom funkcije `sample_gauss_2d`. Pozovite funkciju za učenje te nakon toga provedite klasifikaciju primjera za učenje. Predviđene vjerojatnosti podataka pretvorite u indekse razreda Y (izbjegnite petlju u Pythonu!). Ispišite pokazatelje performanse dobivene pozivima funkcija `eval_perf_binary` te `eval_AP`. Neka vaš ispitni kod izgleda kako slijedi.

```
if __name__ == "__main__":
    np.random.seed(100)

    # get the training dataset
    X, Y_ = data.sample_gauss_2d(2, 100)

    # train the model
    w, b = binlogreg_train(X, Y_)

    # evaluate the model on the training dataset
    probs = binlogreg_classify(X, w, b)
    Y = # TODO

    # report performance
    accuracy, recall, precision = data.eval_perf_binary(Y, Y_)
    AP = data.eval_AP(Y_[probs.argsort()])
    print(accuracy, recall, precision, AP)
```

Ako vaš postupak ne postiže točnost od 100%, pokušajte naći objašnjenje.

Ako je rezultat prihvatljiv, pohranite kod u datoteku `binlogreg.py`. Funkcije `sample_gauss_2d`, `eval_AP` i `eval_perf_binary` pohranite u datoteku `data.py`.

Dodatni zadatci za radoznaće.

- Provjerite poklapaju li se vrijednosti analitičkih gradijenata s njihovim numeričkim aproksimacijama (engl. gradient checking).
- Dodajte regularizacijski gubitak u obliku L2 norme težina w pomnožene hiperparametrom `param_lambda` (parametre b najčešće nema smisla regularizirati). Ne zaboravite izraziti utjecaj te promjene na gradiente.
- Eksperimentirajte s različitim vrijednostima hiperparametara `param_niter`, `param_delta` i `param_lambda`. Pronađite kombinacije hiperparametara za koje vaš program ne uspijeva pronaći zadovoljavajuće rješenje i objasnite što se događa.
- Smislite način za vizualiziranje plohe funkcije gubitka i napredovanje postupka optimizacije.
- Procijenite pristranost (eng. bias) i varijancu vašeg postupka većim brojem eksperimenata na podatcima za testiranje. Podatke za testiranje dobijte uzorkovanjem iste podatkovne distribucije koja je korištена za dobivanje skupa za učenje. Možete li zadati distribuciju podataka za koju će pristranost klasifikatora biti vrlo velika?

3. Grafički prikaz rezultata klasifikacije

Rješenje ove vježbe slobodno preuzmite [ovdje](#).

Napišite funkciju `graph_data` za vizualiziranje rezultata binarne klasifikacije 2D podataka. Neka funkcija ima sljedeće parametre.

```
...
X ... podatci (np.array dimenzija Nx2)
Y_ ... točni indeksi razreda podataka (Nx1)
Y ... predviđeni indeksi razreda podataka (Nx1)
...
```

Upute:

- Podatke iscrtajte funkcijom `plt.scatter`; koristite imenovane argumente `c` i `marker` za zadavanje boje, veličine i oblika simbola koji predstavljaju podatke.
- Točan razred podatka prikažite bojom simbola: razred 0 iscrtajte sivom, a razred 1 bijelom bojom.
- Točno klasificirane podatke ($\text{Y}_\text{==Y}$) prikažite kružićima, a netočno klasificirane podatke kvadratima. Kako biste izbjegli petlju u Pythonu, funkciju `plt.scatter` možete pozvati dva puta: jednom za točno klasificirane, a jednom za netočno klasificirane podatke.

Ispitajte vaš dosadašnji kod na podatcima generiranim dvjema slučajnim razdiobama. Za potrebe prvog testiranja, predviđene indekse razrede podataka (y) odredite proizvoljnom funkcijom odluke, npr:

```
def myDummyDecision(X):
    scores = X[:,0] + X[:,1] - 5
    return scores
```

Vaš novi ispitni program za modul `data` mogao bi izgledati ovako:

```
if __name__=="__main__":
    np.random.seed(100)

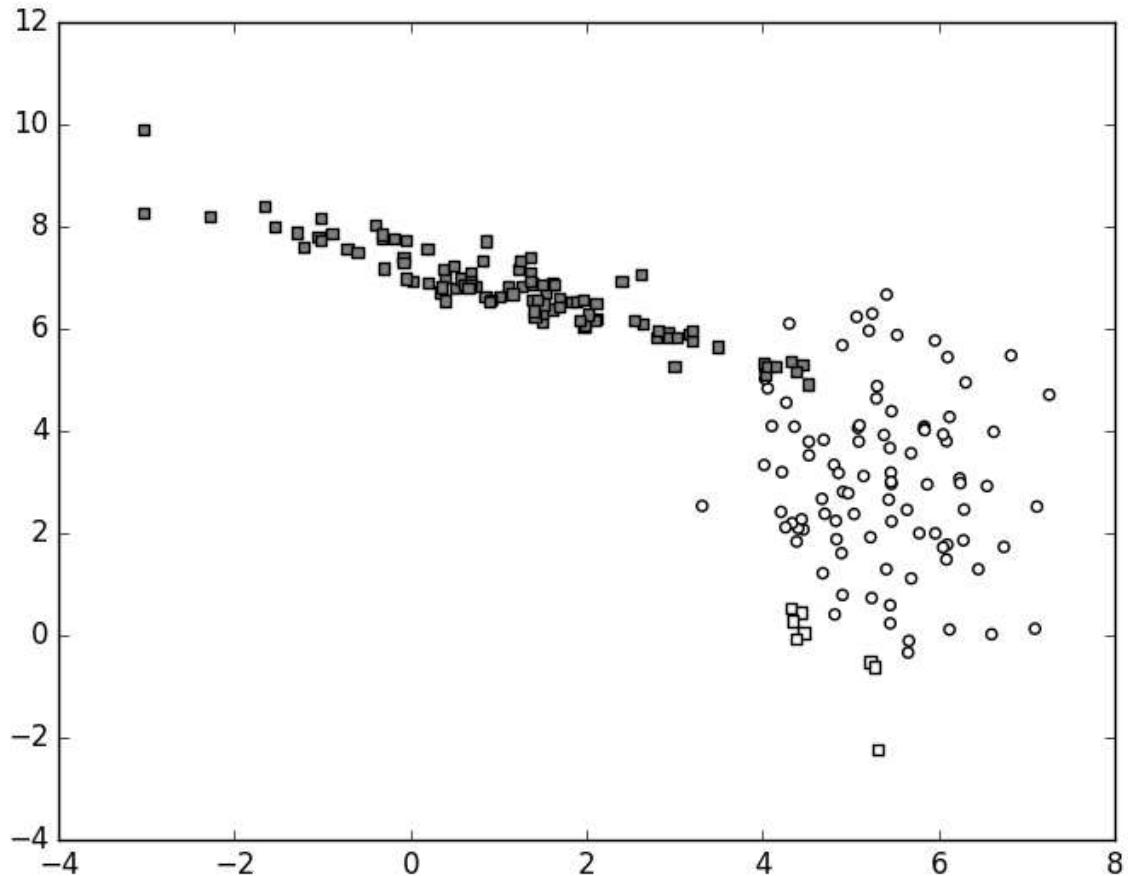
    # get the training dataset
    X,Y_ = data.sample_gauss_2d(2, 100)

    # get the class predictions
    Y = myDummyDecision(X)>0.5

    # graph the data points
    graph_data(X, Y_, Y)

    # show the results
    plt.show()
```

Ako koristimo istu verziju generatora slučajnih brojeva, vaš rezultat sada bi trebao izgledati kao na sljedećoj slici:



Funkcija odluke `myDummyDecision` mijenja vrijednost po dijagonali sjeverozapad-jugoistok, a njena granica prolazi blizu dna nakupine "bijelih". Stoga su "bijeli" podatci uglavnom ispravno klasificirani (kružići), dok su svi "sivi" podatci neispravno klasificirani (kvadrati).

Ako je rezultat ispitivanja prihvatljiv, pohranite kod u datoteku `data.py`.

4. IsCRTavanje funkcije odluke

Rješenje ove vježbe slobodno preuzmite [ovdje](#).

U praksi, funkcije odluke klasifikatora najčešće ne vraćaju predviđeni razred podatka kao što je to bio slučaj kod naše funkcije `myDummyDecision`. Primjerice, funkcija odluke stroja s potpornim vektorima vraća klasifikacijske mjere. Predznak rezultata određuje predviđeni razred, a granicu između dvaju razreda čine podatci u kojima je klasifikacijska mjeru jednaka nuli. Binarna logistička regresija vraća aposteriornu vjerojatnost $P(Y = c_1 | \mathbf{x})$. Ako je ta vjerojatnost veća od 0.5, klasifikator predviđa razred c_1 , a granicu čine podatci u kojima je ta vjerojatnost jednaka 0.5. Jednaku interpretaciju najčešće imaju i funkcije odluke dubokih klasifikatora koji, za razliku od logističke regresije, mogu modelirati i vrlo nelinearne granice među razredima. Kod svih navedenih klasifikatora oblik funkcije odluke pruža kvalitativan uvid u ono što je klasifikator naučio.

Stoga ćemo se u ovoj vježbi posvetiti vizualizaciji funkcije odluke. Naš cilj je razviti potprogram `graph_surface` koji iscrtava plohu zadane skalarne funkcije 2D varijable kako bismo njime mogli iscrtavati plohe funkcije odluke klasifikatora dvodimenzionalnih podataka. Potprogram `graph_surface` treba imati sljedeće argumente:

```
...
fun ... decizijska funkcija (Nx2)->(Nx1)
rect ... željena domena prikaza zadana kao:
       ([x_min,y_min], [x_max,y_max])
offset ... "nulta" vrijednost decizijske funkcije na koju
          je potrebno poravnati središte palete boja;
          tipično imamo:
```

```

offset = 0.5 za probabilističke modele
    (npr. logistička regresija)
offset = 0 za modele koji ne spljošćuju
    klasifikacijske mjere (npr. SVM)
width,height ... rezolucija koordinatne mreže
...

```

Kao i kod potprograma `graph_data` naš kod će biti tanak omotač oko biblioteke `matplotlib`. Kako bismo pozivatelja zaštiti od izvedbenih detalja iscrtavanja, argument `fun` potprograma `graph_surface` predstavlja funkciju odluke (radi se o oblikovnom obrascu Strategija). Funkcija odluke `fun` preslikava podatke u numeričke vrijednosti koje određuju predviđeni razred odgovarajućeg podatka pa će to u praksi biti tanak omotač oko klasifikatora koejg želimo vizualizirati. Kako bismo izbjegli potrebu za višestrukim pozivanjem iz pythonskih petlji, funkcija odluke treba primiti podatkovnu matricu $N \times D$, te vratiti vektor klasifikacijskih mjera dimenzija $N \times 1$. Argument `rect` potprograma `graph_surface` zadaje dio podatkovnog vektorskog prostora u kojem treba vizualizirati funkciju odluke. U praksi, taj argument će ovisiti o rasponu podataka: nema smisla prikazivati funkciju odluke u dijelovima prostora gdje nema podataka za učenje. Potprogram `graph_surface` treba primiti i argument `offset` koji zadaje vrijednost funkcije odluke koja predstavlja granicu između razreda. Za probabilističke modele ta granica će biti 0.5, dok ćemo kod modela koji maksimiziraju marginu tu granicu postaviti na 0. Konačno, potprogram `graph_surface` treba primiti i dimenzije koordinatne mreže koje će definirati zrnatost vizualizacije.

Plohe u `matplotlib` iscrtavamo primjenom funkcija `np.meshgrid` i `plt.pcolormesh`, kao što je to pokazano u ovom [primjeru](#) i službenoj [dokumentaciji](#) biblioteke `Matplotlib`. Kako ta procedura nije baš najintuitivnija, za ovaj zadatak ćemo dati nešto detaljnije instrukcije.

- Izradite 1D koordinatne raspone rezolucije `width×height` u okvirima zadanog pravokutnika `rect` korištenjem funkcije `np.linspace` (ili `np.arange`).
- Izradite 2D koordinatnu mrežu pozivom `np.meshgrid`.
- Izrazite koordinatnu mrežu 2D matricom dimenzija $N \times 2$. Ovaj korak nam treba kako bismo koordinatnu mrežu mogli koristiti kao parametar decizijske funkcije. Koristite metodu `flatten` i funkciju `np.stack`
- Pozovite decizijsku funkciju kako biste dobili njene vrijednosti u točkama koordinatne mreže. Preoblikujte dobivene vrijednosti u oblik kojeg zahtijeva funkcija `plt.pcolormesh` metodom `reshape`.
- IsCRTajte plohu pozivom funkcije `plt.pcolormesh`. Koristite argumente `vmin` i `vmax` tako da paletu centrirate u vrijednosti `offset` te tako da sve vrijednosti funkcije budu u dometu palete.
- IsCRTajte konturu uzduž koje decizijska funkcija poprima vrijednost `offset`. Koristite crnu boju.

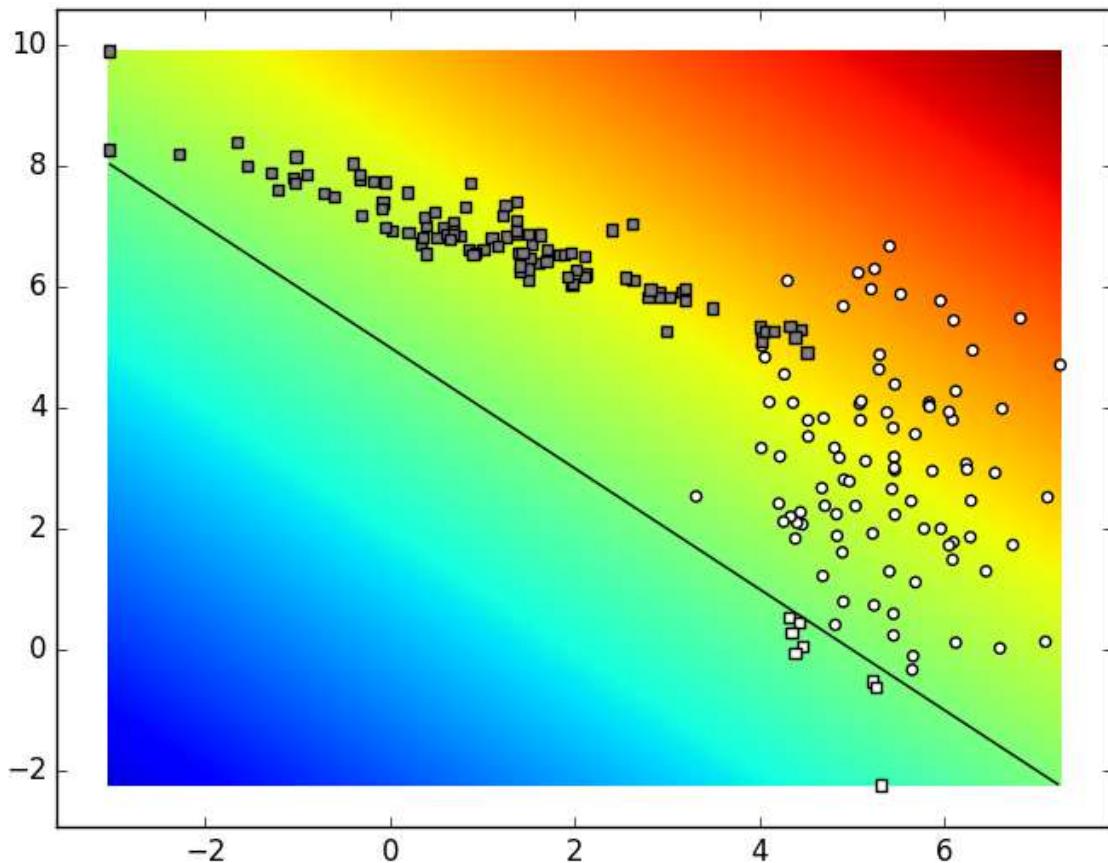
Evo kako bismo predloženim potprogramom iscrtali plohu decizijske funkcije `myDummyDecision` u pravokutniku određenom koordinatama skupa točaka `x`.

```

bbox=(np.min(X, axis=0), np.max(X, axis=0))
graph_surface(myDummyDecision, bbox, offset=0)

```

Ako ovaj odsječak uvrstimo u ispitni program modula `data.py` neposredno prije poziva funkcije `graph_data`, ispod naših podataka osvanut će decizijska ploha u duginim bojama. Točke u kojima decizijska funkcija poprima vrijednost praga označene su crnom linijom. Bijeli podatci iznad crne linije označeni su kružićima (`myDummyClassifier` ih je dobro klasificirao). Bijeli podatci ispod linije i crni podatci iznad linije označeni su kvadratićima (`myDummyClassifier` ih je loše klasificirao).



Ako je rezultat ispitivanja prihvatljiv, pohranite kod u datoteku `data.py`.

5. Grafički prikaz binarne logističke regresije

Sada imamo gotovo sve sastojke za konačnu provjeru naše implementacije binarne logističke regresije. Nedostaje nam samo funkcija koja bi bila prikladna za slanje potprogramu `graph_surface`. U slučaju binarne logističke regresije, željena funkcija trebala bi primiti podatke x , a vratiti vektor aposteriorne vjerojatnosti razreda. Funkcija `binlogreg_classify` izgleda kao obećavajući kandidat, ali ne može se izravno primijeniti jer prima dva argumenta viška (w i b). U programskim jezicima sa statickim funkcijama ovaj izazov bismo teško riješili bez globalnih varijabli. Međutim, u Pythonu ovo možemo lakoćom izvesti primjenom kontekstne funkcije (engl. closure) koju možemo konstruirati pozivom sljedeće funkcije:

```
def binlogreg_decfun(w,b):
    def classify(X):
        return binlogreg_classify(X, w,b)
    return classify
```

Funkcija `binlogreg_decfun` vraća lokalnu funkciju koja pamti kontekst (w, b). Stoga tu lokalnu funkciju možemo koristiti kao argument funkcije `graph_surface` (radi se o oblikovnom obrascu Dekorator). Prikazana fukcionalnost može se sažetije ostvariti lambda izrazom (napravite to za vježbu!). Početnicima savjetujemo da napreduju malim koracima te da lambda izraze koji pamte kontekst počnu koristiti tek nakon što u potpunosti usvoje kontekstne funkcije.

Sada konačno možemo iscrtati rezultate naše izvedbe binarne logističke regresije sljedećim kodom:

```
# instantiate the dataset
# ...
```

```
# train the logistic regression model
# ...

# evaluate the model on the train set
# ...

# recover the predicted classes Y
# ...

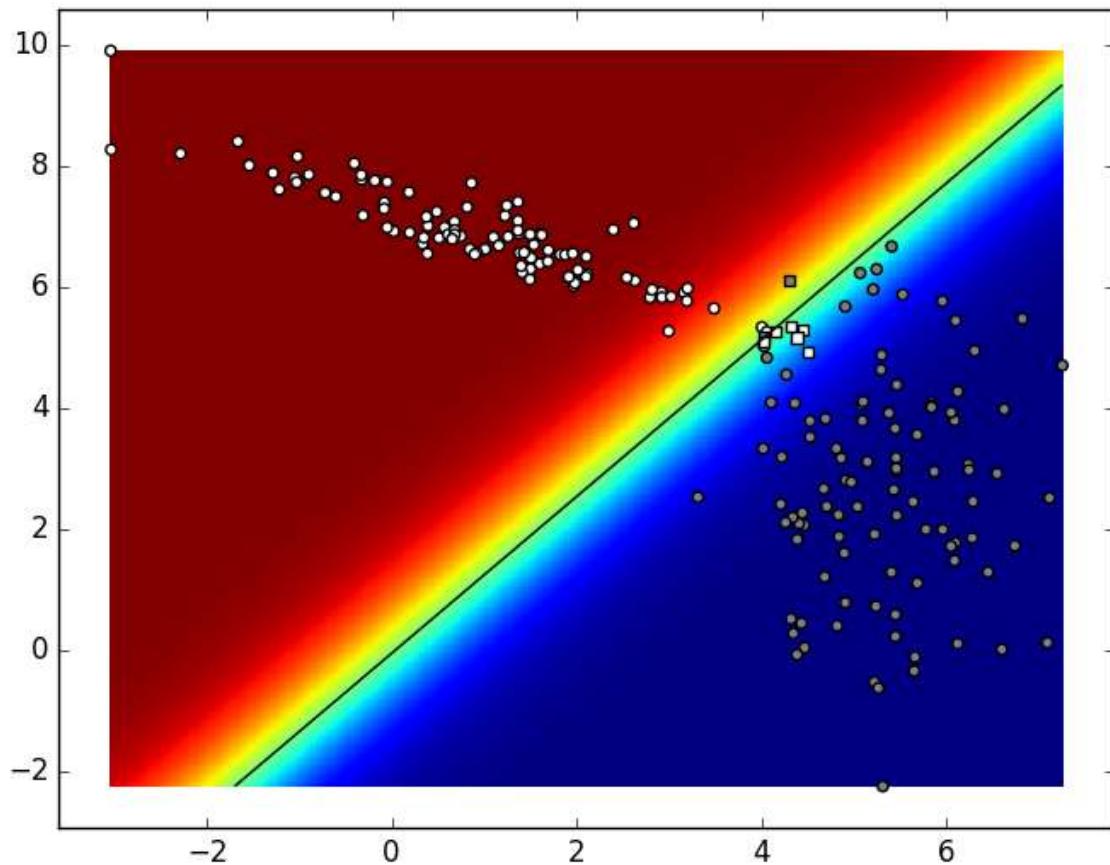
# evaluate and print performance measures
# ...

# graph the decision surface
decfun = binlogreg_decfun(w,b)
bbox=(np.min(X, axis=0), np.max(X, axis=0))
data.graph_surface(decfun, bbox, offset=0.5)

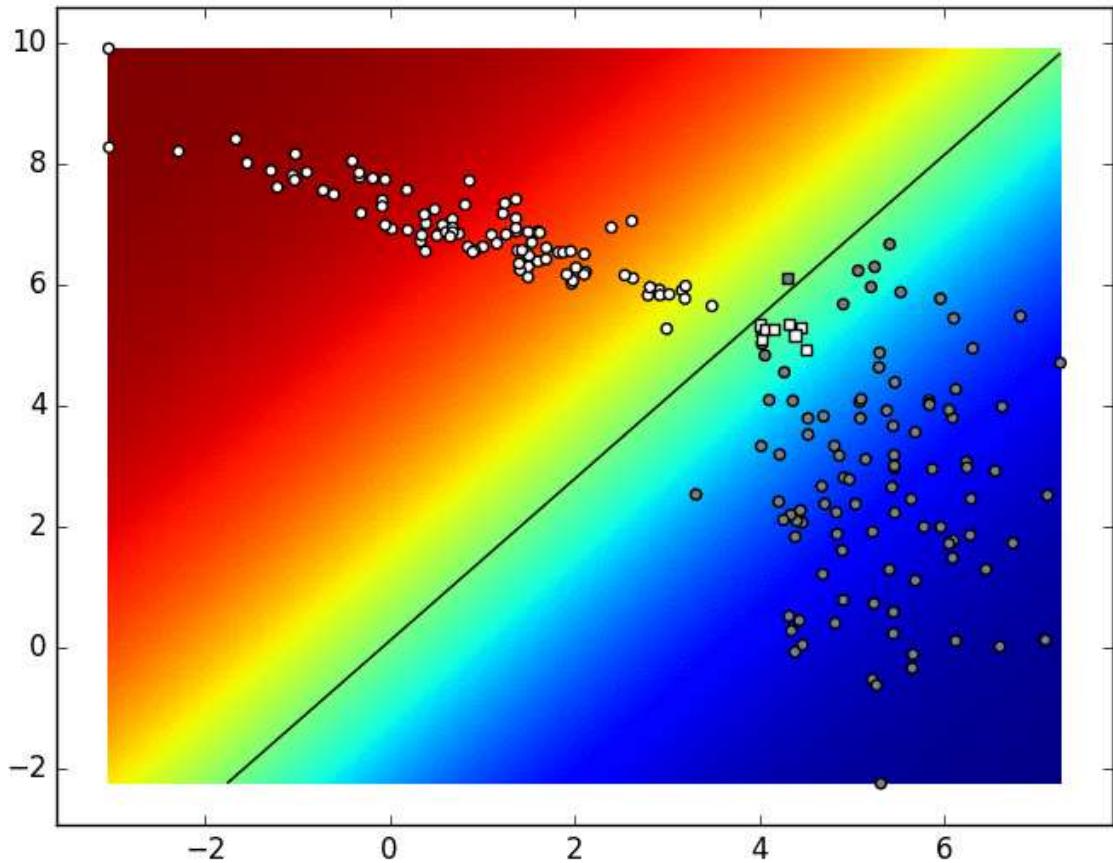
# graph the data points
data.graph_data(X, Y_, Y, special=[])

# show the plot
plt.show()
```

Ovisno o stanju generatora slučajnih brojeva te vrijednostima hiperparametara vaš rezultat mogao bi izgledati kao na sljedećoj slici:



Ako sliku izradimo u svakoj iteraciji postupka, dobivamo sljedeću animaciju. Pokušajte objasniti razlike između konačnog stanja ove animacije i funkcije odluke koja je prikazana na prethodnoj slici.



6. Višerazredna logistička regresija

Naš konačni zadatak je napisati kod za učenje i primjenu logističke regresije nad podatcima proizvoljnog broja razreda. Potrebno je napisati funkciju `logreg_train` koja kao argumente prima podatkovnu matricu `x` i matricu indeksa točnih razreda `y_`. Na izlazu je potrebno vratiti parametre logističke regresije `w` i `b`. Dimenzije tih parametara trebaju ovisiti o broju razreda `c` kojeg možete odrediti izrazom `max(Y_) + 1`.

Struktura tijela petlje gradijentnog spusta trebala bi izgledati ovako:

```
# eksponencirane klasifikacijske mjere
# pri računanju softmaxa obratite pažnju
# na odjeljak 4.1 udžbenika
# (Deep Learning, Goodfellow et al)!
scores = ...      # N x C
expscores = ...  # N x C

# nazivnik softmaxa
sumexp = ...     # N x 1

# logaritmirane vjerojatnosti razreda
probs = ...      # N x C
logprobs = ...   # N x C

# gubitak
loss = ...        # scalar

# dijagnostički ispis
if i % 10 == 0:
    print("iteration {}: loss {}".format(i, loss))

# derivacije komponenata gubitka po mjerama
dL_dz = ...       # N x C

# gradijenti parametara
```

```

grad_W = ...      # C x D (ili D x C)
grad_b = ...      # C x 1 (ili 1 x C)

# poboljšani parametri
W += -param_delta * grad_W
b += -param_delta * grad_b

```

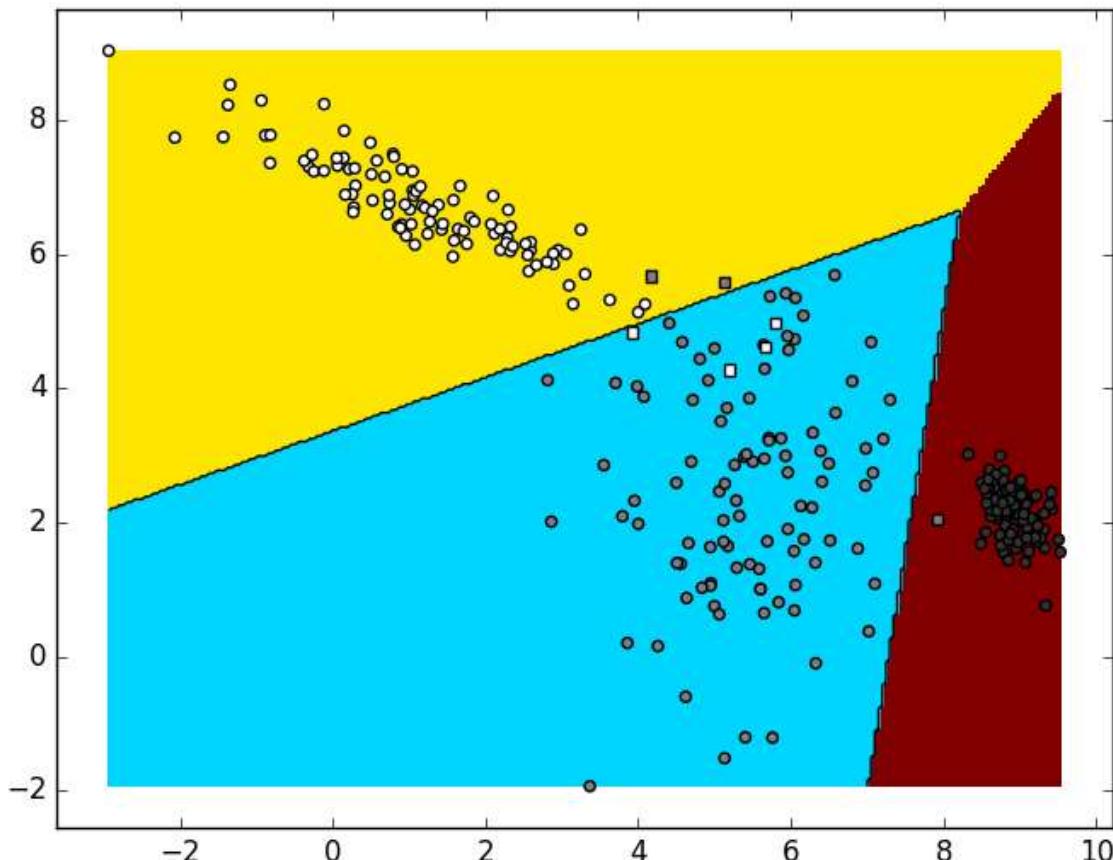
Napišite funkciju `logreg_classify` koja kao rezultat vraća matricu dimenzija $N \times C$ gdje svaki redak i sadrži vjerojatnosti klasificiranja podatka \mathbf{x}_i u razrede c_j , $j \in \{0, 1, \dots, C - 1\}$.

Napišite funkciju `eval_perf_multi(Y, Y_)` koja na temelju predviđenih i točnih indeksa razreda određuje pokazatelje klasifikacijske performanse: ukupnu točnost klasifikacije, matricu **zabune** (engl. confusion matrix) u kojoj retci odgovaraju točnim razredima a stupci predikcijama te vektore preciznosti i odziva **pojedinih** razreda. U implementaciji prvo izračunajte matricu zabune, a onda sve ostale pokazatelje na temelju nje.

Konstruirajte dekorator oko `logreg_classify` koji omogućava evaluiranje modela iz funkcije `graph_surface`. Decizijsku plohu možete prikazati na više načina:

- kao vjerojatnost nekog određenog razreda,
- kao maksimalnu vjerojatnost klasificiranja u bilo koji razred,
- kao razred koji ima najveću aposteriornu vjerojatnost.

Ispobrijte vaš kod najprije za slučaj $C=2$. Trebali biste dobiti vrlo slične rezultate kao i u binarnom slučaju. Nakon toga, provjerite što se zbiva kad vašem kodu pošaljete podatke uzorkovane iz tri slučajne Gaussove distribucije. Obratite pažnju da, kao i ranije, podatke iz svake pojedine distribucije valja označiti zasebnim razredom. Ovisno o stanju generatora slučajnih brojeva te vrijednostima hiperparametara vaš rezultat mogao bi izgledati kao na sljedećoj slici:



Ako je rezultat ispitivanja prihvatljiv, pohranite kod u datoteku `logreg.py`.

Ove stranice su rezultat rada na istraživačkom projektu [MULTICLOD](#) (I-2433-2014) Hrvatske zaklade za znanost.

Stranice su izrađene [vi-jem](#) i [geditom](#).

Posljednja promjena: Wednesday, 18-Oct-2017 11:25:41 CEST

Svi komentari su dobrodošli: sinisa.segvic@fer.hr

[Povratak](#)