

## 6. Logistička regresija

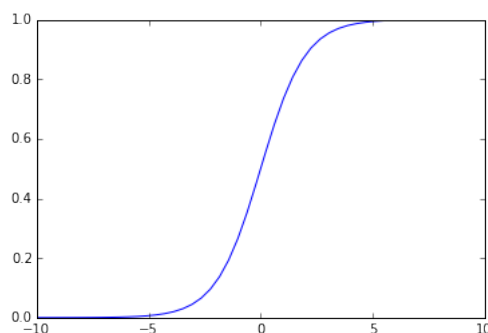
Strojno učenje 1, UNIZG FER, ak. god. 2021./2022.

Jan Šnajder, natuknice s predavanja, v1.5

### 1 Model logističke regresije

- **Logistička (sigmoidalna) funkcija:**

$$\sigma(\alpha) = \frac{1}{1 + \exp(-\alpha)}$$



- Funkcija je derivabilna:

$$\frac{\partial \sigma(\alpha)}{\partial \alpha} = \frac{\partial}{\partial \alpha} (1 + \exp(-\alpha))^{-1} = \sigma(\alpha)(1 - \sigma(\alpha))$$

- Model logističke regresije:

$$h(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \phi(\mathbf{x})) = \frac{1}{1 + \exp(-\mathbf{w}^T \phi(\mathbf{x}))} = P(y = 1 | \mathbf{x})$$

⇒ izlaz modela možemo tumačiti kao vjerojatnost da primjer pripada klasi  $y = 1$

- Ovo je primjer **poopćenog linearnog modela** (*generalized linear model, GLM*)
- GLM – linearni modeli s (nelinarnom) **aktivacijskom funkcijom**  $f$ :

$$h(\mathbf{x}; \mathbf{w}) = f(\mathbf{w}^T \phi(\mathbf{x}))$$

gdje je  $f : \mathbb{R} \rightarrow [0, 1]$  ili  $f : \mathbb{R} \rightarrow (0, 1)$  ili  $f : \mathbb{R} \rightarrow [-1, +1]$  ili  $f : \mathbb{R} \rightarrow (-1, +1)$

## 2 Pogreška unakrsne entropije

- Izlaz modela je **Bernoullijeva varijabla**:

$$P(y|\mu) = \begin{cases} \mu & \text{ako } y = 1 \\ 1 - \mu & \text{inače} \end{cases} = \mu^y(1 - \mu)^{1-y}$$

- U našem slučaju,  $y$  je oznaka primjera, a  $\mu$  je izlaz modela, tj.  $\mu = h(\mathbf{x}; \mathbf{w})$ , pa:

$$P(y^{(i)}|\mathbf{x}^{(i)}) = h(\mathbf{x}; \mathbf{w})^y(1 - h(\mathbf{x}; \mathbf{w}))^{1-y}$$

- Log-izglednost oznaka iz skupa označenih primjera:

$$\begin{aligned} \ln P(\mathbf{y}|\mathbf{X}, \mathbf{w}) &= \ln \prod_{i=1}^N p(y^{(i)}|\mathbf{x}^{(i)}) = \\ &= \sum_{i=1}^N \left( y^{(i)} \ln h(\mathbf{x}^{(i)}; \mathbf{w}) + (1 - y^{(i)}) \ln (1 - h(\mathbf{x}^{(i)}; \mathbf{w})) \right) \end{aligned}$$

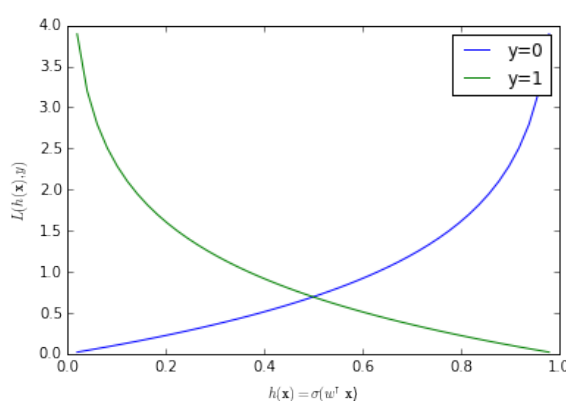
- Empirijska pogreška je negativna log-izglednost:

$$E(\mathbf{w}|\mathcal{D}) = \frac{1}{N} \sum_{i=1}^N \left( -y^{(i)} \ln h(\mathbf{x}^{(i)}; \mathbf{w}) - (1 - y^{(i)}) \ln (1 - h(\mathbf{x}^{(i)}; \mathbf{w})) \right)$$

⇒ **pogreška unakrsne entropije (cross-entropy error)**

- **Gubitak unakrsne entropije (cross-entropy loss):**

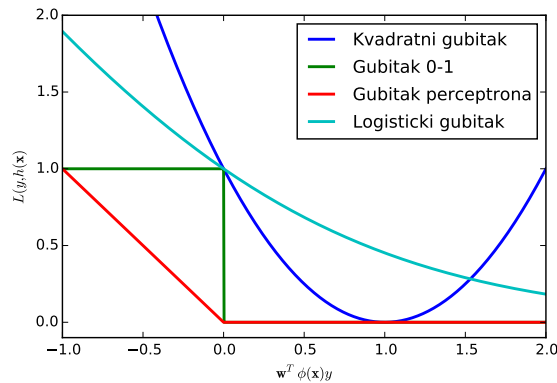
$$L(y, h(\mathbf{x})) = -y \ln h(\mathbf{x}) - (1 - y) \ln (1 - h(\mathbf{x}))$$



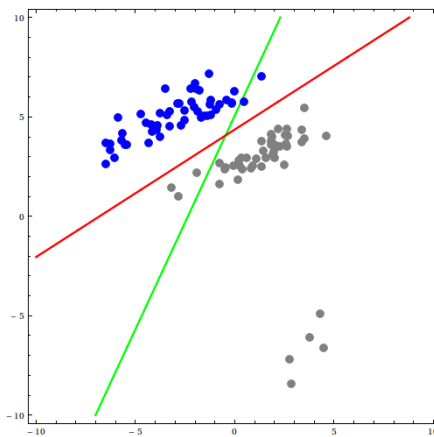
- Reformulacija  $y \in \{0, 1\} \rightarrow y \in \{-1, +1\}$  i skaliranje sa  $1/\ln 2$ :

$$L(y, h(\mathbf{x})) = \frac{1}{\ln 2} \ln (1 + \exp(-y\mathbf{w}^T \phi(\mathbf{x})))$$

- Usporedba funkcija gubitaka:



- Logistička regresija robusnija je od modela linearne regresije:



- Minimizacija u zatvorenoj formi nije moguća  $\Rightarrow$  iterativna optimizacija

### 3 Gradijentni spust

- **Gradijentni spust** – minimum nalazimo krećući se u smjeru suprotnom od gradijenta:

$$\mathbf{x} \leftarrow \mathbf{x} - \eta \nabla f(\mathbf{x})$$

- $\eta$  je **stopa učenja**: prevelika  $\eta \Rightarrow$  divergencija; premalena  $\eta \Rightarrow$  spora konvergencija
- Želimo **globalnu konvergenciju** (konvergencija uvijek i svugdje)
- Ostvarivo **linijskim pretraživanjem** –  $\eta$  koji minimizira  $f(\mathbf{x})$  u smjeru spusta  $\Delta\mathbf{x}$ :

$$g(\eta) = f(\mathbf{x} + \eta\Delta\mathbf{x})$$

- Pronađeni optimum bit će globalni optimum ako je  $f(\mathbf{x})$  **konveksna**
- Funkcija  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  je **konveksna** akko

(1) Njezina domena  $\text{dom}(f)$  je **konveksni skup**:

Za svaki  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \text{dom}(f)$  i za svaki  $\alpha_1, \dots, \alpha_n$  takav da  $\sum_i \alpha_i = 1$  vrijedi:

$$\sum_{i=1}^n \alpha_i \mathbf{x}_i \in \text{dom}(f)$$

(2) Za svaki  $\mathbf{x}_1, \mathbf{x}_2 \in \text{dom}(f)$  i svaki  $\alpha \in [0, 1]$  vrijedi:

$$f(\mathbf{x}) = f(\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2) \leq \alpha f(\mathbf{x}_1) + (1 - \alpha) f(\mathbf{x}_2)$$

• Empirijska pogreška je konveksna  $\Leftrightarrow$  funkcija gubitka  $L$  je konveksna

• Dvije varijante gradijentnog spusta:

– **Batch** (grupni):  $\mathbf{w} \leftarrow \mathbf{w} - \eta \sum_{i=1}^N \nabla L(y^{(i)}, h(\mathbf{x}^{(i)}; \mathbf{w}))$

– **Stohastički (SGD)**:  $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla L(y^{(i)}, h(\mathbf{x}^{(i)}; \mathbf{w}))$

• SGD je pogodan za on-line učenje (big data, data streams)

## 4 Gradijentni spust za logističku regresiju

• Gradijent funkcije gubitka i funkcije pogreške:

$$E(\mathbf{w}|\mathcal{D}) = \frac{1}{N} \sum_{i=1}^N \left( -y^{(i)} \ln h(\mathbf{x}^{(i)}; \mathbf{w}) - (1 - y^{(i)}) \ln (1 - h(\mathbf{x}^{(i)}; \mathbf{w})) \right)$$

$$\nabla_{\mathbf{w}} E(\mathbf{w}|\mathcal{D}) = \frac{1}{N} \sum_{i=1}^N \nabla L(y^{(i)}, h(\mathbf{x}^{(i)}; \mathbf{w}))$$

$$\nabla L(y, h(\mathbf{x})) = \left( -\frac{y}{h(\mathbf{x})} + \frac{1-y}{1-h(\mathbf{x})} \right) h(\mathbf{x})(1-h(\mathbf{x})) \phi(\mathbf{x}) = (h(\mathbf{x}) - y) \phi(\mathbf{x})$$

$$\nabla E(\mathbf{w}|\mathcal{D}) = \sum_{i=1}^N (h(\mathbf{x}^{(i)}) - y^{(i)}) \phi(\mathbf{x}^{(i)})$$

(faktor  $1/N$  može se apsorbirati u stopu učenja  $\eta$ )

### Logistička regresija (grupni gradijentni spust)

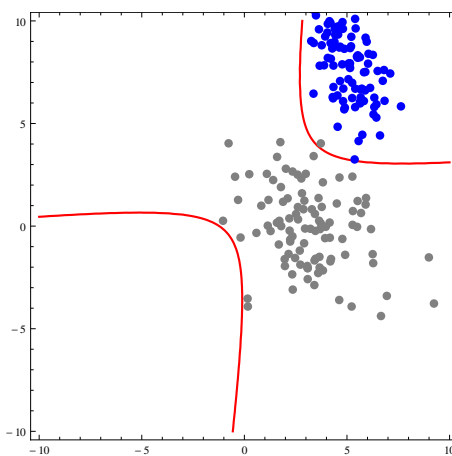
- 1:  $\mathbf{w} \leftarrow (0, 0, \dots, 0)$
- 2: **ponavljaj** do konvergencije
- 3:  $\Delta \mathbf{w} \leftarrow (0, 0, \dots, 0)$
- 4: **za**  $i = 1, \dots, N$
- 5:  $h \leftarrow \sigma(\mathbf{w}^T \phi(\mathbf{x}^{(i)}))$
- 6:  $\Delta \mathbf{w} \leftarrow \Delta \mathbf{w} - (h - y^{(i)}) \phi(\mathbf{x}^{(i)})$
- 7:  $\eta \leftarrow$  optimum linijskim pretraživanjem u smjeru spusta  $\Delta \mathbf{w}$
- 8:  $\mathbf{w} \leftarrow \mathbf{w} + \eta \Delta \mathbf{w}$

## Logistička regresija (stohastički gradijentni spust)

- 1:  $\mathbf{w} \leftarrow (0, 0, \dots, 0)$
- 2: **ponavljaj** do konvergencije
- 3: slučajno permutiraj primjere u  $\mathcal{D}$
- 4: **za**  $i = 1, \dots, N$
- 5:  $h \leftarrow \sigma(\mathbf{w}^T \phi(\mathbf{x}^{(i)}))$
- 6:  $\Delta \mathbf{w} \leftarrow -(h - y^{(i)}) \phi(\mathbf{x}^{(i)})$
- 7:  $\eta \leftarrow$  optimum linijskim pretraživanjem u smjeru spusta  $\Delta \mathbf{w}$
- 8:  $\mathbf{w} \leftarrow \mathbf{w} + \eta \Delta \mathbf{w}$

## 5 Regularizirana regresija

- Prednosti regularizacije:
  - Sprječavanje pretjerane nelinearnosti
  - Suzbijanje nepotrebnih značajki
  - Sprječavanje otvrdnjavanja sigmoide kod linearno odvojivih problema
- Primjer prenaučivosti ( $n = 2$ ,  $\phi(\mathbf{x}) = (1, x_1, x_2, x_1x_2, x_1^2, x_2^2)$ ):



- L2-regularizirana pogreška:

$$E_R(\mathbf{w}|\mathcal{D}) = \sum_{i=1}^N \left( -y^{(i)} \ln h(\mathbf{x}^{(i)}) - (1 - y^{(i)}) \ln (1 - h(\mathbf{x}^{(i)})) \right) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- Ažuriranje težina:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \left( \sum_{i=1}^N (h(\mathbf{x}^{(i)}) - y^{(i)}) \phi(\mathbf{x}^{(i)}) + \lambda \mathbf{w} \right)$$

ekvivalentno:

$$\mathbf{w} \leftarrow \mathbf{w}(1 - \eta\lambda) - \eta \sum_{i=1}^N (h(\mathbf{x}^{(i)}) - y^{(i)}) \phi(\mathbf{x}^{(i)})$$

- Napomena: Težina  $w_0$  se ne regularizira

### L2-regularizirana logistička regresija (grupni gradijentni spust)

```

1:  $\tilde{\mathbf{w}} \leftarrow (0, 0, \dots, 0)$  //  $\tilde{\mathbf{w}}$  je prošireni vektor  $(w_0, \mathbf{w})$ 
2: ponavljaj do konvergencije
3:    $\Delta w_0 \leftarrow 0$ 
4:    $\Delta \mathbf{w} \leftarrow (0, 0, \dots, 0)$ 
5:   za  $i = 1, \dots, N$ 
6:      $h \leftarrow \sigma(\tilde{\mathbf{w}}^T \phi(\tilde{\mathbf{x}}^{(i)}))$ 
7:      $\Delta w_0 \leftarrow \Delta w_0 - (h - y^{(i)})$ 
8:      $\Delta \mathbf{w} \leftarrow \Delta \mathbf{w} - (h - y^{(i)}) \phi(\mathbf{x}^{(i)})$ 
9:    $\eta \leftarrow$  optimum linijskim pretraživanjem u smjeru spusta  $\Delta \tilde{\mathbf{w}}$ 
10:   $w_0 \leftarrow w_0 + \eta \Delta w_0$ 
11:   $\mathbf{w} \leftarrow \mathbf{w}(1 - \eta\lambda) + \eta \Delta \mathbf{w}$ 

```

### L2-regularizirana logistička regresija (stohastički gradijentni spust)

```

1:  $\tilde{\mathbf{w}} \leftarrow (0, 0, \dots, 0)$  //  $\tilde{\mathbf{w}}$  je prošireni vektor  $(w_0, \mathbf{w})$ 
2: ponavljaj do konvergencije:
3:   slučajno permutiraj primjere u  $\mathcal{D}$ 
4:   za  $i = 1, \dots, N$ 
5:      $h \leftarrow \sigma(\tilde{\mathbf{w}}^T \phi(\tilde{\mathbf{x}}^{(i)}))$ 
6:      $\Delta w_0 \leftarrow -(h - y^{(i)})$ 
7:      $\Delta \mathbf{w} \leftarrow -(h - y^{(i)}) \phi(\mathbf{x}^{(i)})$ 
8:    $\eta \leftarrow$  optimum linijskim pretraživanjem u smjeru spusta  $\Delta \tilde{\mathbf{w}}$ 
9:    $w_0 \leftarrow w_0 + \eta \Delta w_0$ 
10:   $\mathbf{w} \leftarrow \mathbf{w}(1 - \eta\lambda) + \eta \Delta \mathbf{w}$ 

```