Statistical Machine Learning – Exercise 2 – Michael Meidlinger

## 1 Bayes Classifier

a) We have $\quad p(x,y) = p(y|x)p(x) = p(x|y)P(y) \qquad (1)$

- $c_1$ is NOT equivalent. In fact, $p(x)$ is independent of $y$ so that $\arg\max_y p(x)$ can be defined to output any arbitrary number

- $c_2$ is NOT equivalent. The output is independent of $x$ always the least likely label

- $c_3$ is equivalent because of (1) and the fact that $p(x)$ is positive and independent of $y$

- $c_4$ is NOT equivalent, since $p(x|y) = p(y|x)\frac{p(x)}{P(y)}$

b.) Using a similar argumentation as above we have

$\quad c_5, c_7, c_8, c_9, c_{10}, c_{12} \quad$ equivalent, the rest not

## 2 Gaussian Discriminant Analysis

a) $\quad c(x) = \arg\max_{y \in \mathcal{Y}} p(y|x) = \arg\max_{y \in \mathcal{Y}} p(x|y)P(y) = \left|\, y \in \{-1,1\} \,\right| = \text{sign}\left( \log \frac{p(x|1)p(y=1)}{p(x|-1)p(y=-1)} \right)$

$\quad = \text{sign}\left( \log\left( \exp\left( -\tfrac{1}{2}(x-\mu_1)^T \Sigma^{-1}(x-\mu_1) + \tfrac{1}{2}(x-\mu_{-1})^T \Sigma^{-1}(x-\mu_{-1}) + \log\frac{p(y=1)}{p(y=-1)} \right) \right) \right)$

quadratic $x$ terms cancel, quadratic $\mu$ terms can be dropped (independent of $y$)

$\quad = \text{sign}\left( \underbrace{x^T \Sigma^{-1}\mu_1 - x^T \Sigma^{-1}\mu_{-1}}_{x^T\left(\underbrace{\Sigma^{-1}\mu_1 - \Sigma^{-1}\mu_{-1}}_{:=w}\right)} + \underbrace{\log\frac{p(y=1)}{p(y=-1)}}_{\theta} \right)$

$\quad = \text{sign}\left( \langle x, w \rangle + \theta \right) \qquad\qquad q.e.d.$

b.) With a few examples, the exact distributions cannot be estimated precisely. Thus one has to resort to parametric models of the distribution.

# 3 Robustness of the Perceptron

The robustness $\rho$ of a classifier $g$ (with respect to $\mathcal{D}$) is the largest amount by which we can perturb the training samples without changing the predictions of $g$:

$$g(x^i + \varepsilon) = g(x^i) \quad \forall i \; \forall \varepsilon: \|\varepsilon\| < \rho$$



For a linear classifier $\rho$ is the smallest distance of any $x^i \in D_x$ from the decision boundary. For the specified training data with points $x^i \in \mathbb{R}^2$ we can distinguish the following cases:

a) $\binom{a}{b}$ is an element of the ray $r \Rightarrow D$ is not linearly separable and the perceptron algorithm won't converge

b) $\binom{a}{b}$ is an element of $R_2 := \left\{ x \in \mathbb{R}^2 \; \middle| \; \|x - a\| < \frac{\|b - a\|}{2} \wedge \|x - b\| < \frac{\|b - a\|}{2} \right\}$

For that case, $\rho$ is given by

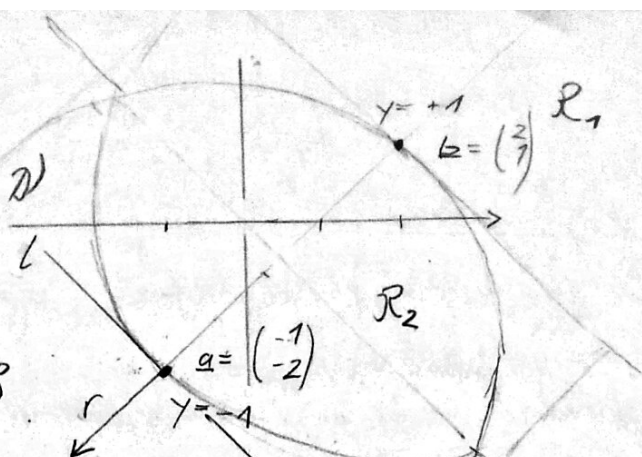$$\rho = \frac{\|x - a\|}{2} = \frac{1}{2}\left\|\binom{a+1}{b+2}\right\| = \frac{1}{2}\sqrt{(a+1)^2 + (b+2)^2}$$

c) $\binom{a}{b}$ is in $\Rightarrow R_1$

$$\rho = \frac{1}{2}\frac{\|b - a\|}{2} = \frac{1}{2}\sqrt{\underbrace{3^2 + 3^2}_{2.9}} = \frac{3}{\sqrt{2}}$$

d) $\binom{a}{b} \in R_2 \setminus \{R_1 \cup R_2 \cup r\}$

The decision boundary is the line $d: \frac{x_0 + x}{2} + \lambda(b - x)$ where

$$x_0 = \arg\min_{y \in \{a, b\}} \|y - x\| \qquad \Rightarrow \quad \rho = \frac{\left(\frac{x - x_0}{2}\right) \times (b - x)}{\|b - x\|} = \left\{$$

- $\rho_{max}$ is thus given by $\frac{3}{\sqrt{2}}$

# 4 Perception Training as Convex Optimization

Comparing Algorithm 1 with "Stochastic Gradient Descent" (ml 2014_09.pdf)
we identify:

$$-\eta \nabla f_i(w_t) \text{ with random } i$$

$$w_{t+1} \leftarrow w_t - \eta_6 \vee \quad \Longleftrightarrow \quad w_{t+1} \leftarrow w_t + y\,x$$

$$\Longleftrightarrow \quad -\eta_t n \nabla f_i(w_t) = \begin{cases} y^i x^i, & y\langle w_t, x\rangle \leq 0 \\ 0, & \text{otherwise} \end{cases}$$

$$[x]^+ = \begin{cases} x, & x > 0 \\ 0, & \text{otherwise} \end{cases}$$

Average of inner products is the cost function to be minimized

$$\Rightarrow \eta_t = \text{const} = 1$$

$$f_i(w) = \frac{1}{n}\left[-y^i \langle w, x^i\rangle\right]^+ \quad \Rightarrow \quad f(w) = \sum_{i=1}^{n} f_i(w) = \overbrace{\frac{1}{n}\sum_{i=1}^{n}\left[-y^i\langle w, x^i\rangle\right]^+}$$

Advantages:

Shortcomings:

# 5 Hard-Margin SVM Dual

$$\min_{w \in \mathbb{R}^d,\, b \in \mathbb{R}} \tfrac{1}{2}\|w\|^2 \quad \text{subject to} \quad y^i\left(\langle w, x^i\rangle + b\right) \geq 1$$

We compute the Lagrangian:

$$\overbrace{\sum_i \alpha_i - \sum_i \alpha_i y^i \langle w, x^i\rangle - b\sum_i \alpha_i y^i}$$

$$\mathcal{L}(w, b, \alpha) = \tfrac{1}{2}\|w\|^2 + \sum_i \alpha_i\left(1 - y^i\left(\langle w, x^i\rangle + b\right)\right), \quad h(\alpha) = \min_{(w,b)} \mathcal{L}(w, b, \alpha)$$

$$\Rightarrow \quad 0 = \frac{\partial}{\partial w}\mathcal{L} = w - \sum_i \alpha_i y_i x^i \quad \Rightarrow \quad w = \sum_i \alpha_i y^i x^i$$

$$0 = \frac{\partial}{\partial b}\mathcal{L} = \sum_i \alpha_i y^i$$

↳ Insert back:

$$\left\|\sum_i \alpha_i y^i x^i\right\|^2$$

$$h(\alpha) = \tfrac{1}{2}\left\|\sum_i \alpha_i y^i x^i\right\|^2 + \sum_i \alpha_i - \overbrace{\sum_i \alpha_i y^i \left\langle \sum_j \alpha_j y^j x^j, x^i\right\rangle}$$

$$= -\tfrac{1}{2}\left\|\sum_i \alpha_i y^i x^i\right\|^2 + \sum_i \alpha_i$$

$$= -\tfrac{1}{2}\left\langle \sum_i \alpha_i y^i x^i, \sum_j \alpha_j y^j x^j\right\rangle + \sum_i \alpha_i = -\tfrac{1}{2}\sum_i \sum_j \alpha_i \alpha_j y^i y^j \langle x^i, x^j\rangle + \sum_i \alpha_i$$

Dual Problem:

$$\min_{\alpha > 0} h(\alpha) \quad \text{subject to}$$

$$\sum_i \alpha_i y^i = 0,$$

# 6 Missing Proofs

$$= \max \{ f_1(w), \ldots, f_K(w) \}$$

a) We need to show $f_K(w) \geqslant f_K(w_0) + \langle v, w-w_0 \rangle \implies f(w) \geqslant f_K(w_0) + \langle v, w-w_0 \rangle \quad \forall k$

$\searrow$ Use assumption

$$f(w) = \max_{k=1,\ldots,K} f_{K'}(w) \geqslant f_K(w) \geqslant \underbrace{f_K(w_0)}_{= f(w_0) \ \text{(assumption)}} + \langle v, w-w_0 \rangle = f(w_0) + \langle v, w-w_0 \rangle \qquad \square$$

b)

$$\overset{w_t - \eta v}{\overset{|}{\| w_{t+1} - w^* \|}}^2 = \| w_{t+1} \|^2 + \| w^* \|^2 - 2 \langle w_{t+1}, w^* \rangle = \| w_t - w^* \|^2$$

$$= \| w_t \|^2 + \| w^* \|^2 - 2 \langle w_t, w^* \rangle + 2\eta \langle v, w^* \rangle + \eta^2 \| v \|^2 \qquad (1)$$

$$(2)$$

$$\| w_t - w^* \|^2 = \| w_t \|^2 + \| w^* \|^2 - 2 \langle w_t, w^* \rangle$$

$v$ is a subgradient at $w_t$ $\iff$ $f(w) \geqslant f(w_t) + \langle v, w - w_t \rangle$, in particular, for $w = w_{t+1}$

$$f(w_{t+1}) - f(w_t) \geqslant \langle v, -\eta v \rangle \iff f(w_t) - f(w_{t+1}) \leq \eta \| v \|^2 \qquad (3)$$

$$f(\theta w_{t+1} + (1-\theta) w_t) = f(-\eta \theta v + w_t)$$

$$\leq \theta f(w_{t+1}) + (1-\theta) f(w_t) = f(w_t) + \theta [ f(w_{t+1}) - f(w_t) ]$$

We now show $\| w_{t+1} - w^* \| < \| w_t - w^* \|$ by showing $\| w_t - w^* \|^2 - \| w_{t+1} - w^* \|^2 > 0$, which is equivalent to $(1),(2)$ showing

$$\| v \|^2 \eta < 2 \langle v, w^* \rangle$$