# 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 $\arg\max\limits_y p(x)$ 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

$\qquad c_7, c_8, c_9, c_{10}, c_{12} \quad$ equivalent, the rest not
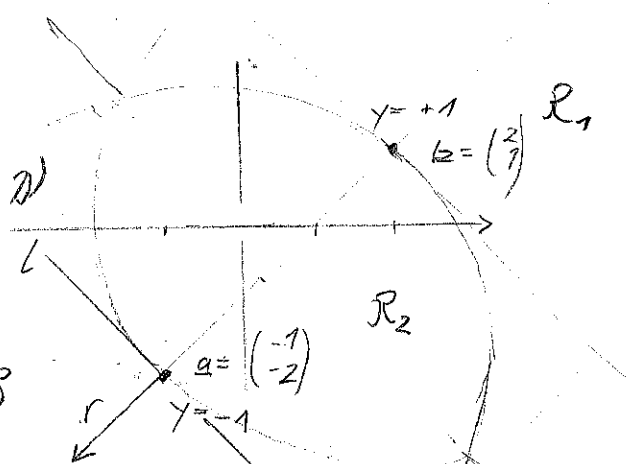
## 2 Gaussian Discriminant Analysis

a) $\quad \angle(x) = \arg\max\limits_{y \in \mathcal{Y}} p(y|x) = \arg\max\limits_{y \in \mathcal{Y}} p(x|y)p(y) = \Big| y \in \{-1,1\} \Big| = \mathrm{sign}\left( \log \frac{p(x|1)p(y=1)}{p(x|-1)p(y=-1)} \right)$

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

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

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

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

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

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



For a linear classifier, $\varrho$ is the smallest distance of any $x^i \in D_x$ to 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 $R_2 := \left\{ x \in \mathbb{R}^2 \mid \|x - a\| < \frac{\|b-a\|}{2} \wedge \|x-b\| < \frac{\|b-a\|}{2} \right\}$

For that case, $\varrho$ is given by

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

c) $\binom{a}{b}$ is in $R_1 \Rightarrow \varrho = \frac{1}{2} \frac{\|b-a\|}{2} = \frac{1}{2}\sqrt{\underbrace{3^2 + 3^2}_{2 \cdot 9}} = \frac{3}{\sqrt{2}}$

d) $\binom{a}{b} \in \mathbb{R}^2 \setminus \{R_1 \cup R_2 \cup r\}$

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

$$x_0 = \underset{y \in [a,b]}{\arg \min} \|y - x\| \qquad \Rightarrow \varrho = \frac{\left(\frac{x - x_0}{2}\right) \times (b - x)}{\|b - x\|} = \Big\{$$

• $\varrho_{max}$ is thus given by $\frac{3}{\sqrt{2}}$

## 4 Perceptron Training as Convex Optimization

Comparing Algorithm 1 with "Stochastic Gradient Descent" (ml2014_09.pdf) we identify:

$$w_{t+1} \leftarrow w_t - \eta_t \underset{\substack{\uparrow \\ -n\nabla f_i(w_t) \text{ with random } i}}{\nabla} \iff w_{t+1} \leftarrow w_t + y\,x$$

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

$$\Rightarrow \eta_t = \text{const} = 1 \qquad , \quad [x]^+ = \begin{cases} x, & x > 0 \\ 0, & \text{otherwise} \end{cases}$$

Average of inner products is the cost function to be minimized

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

Advantages :

Shortcomings :

## 5 Hard-Margin SVM Dual

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

We compute the Lagrangian:
$$\overbrace{\sum_i \alpha_i - \sum_i \alpha_i y^i\langle w, x^i\rangle - b\sum_i \alpha_i y^i}$$

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

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

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

$\hookrightarrow$ Insert back:

$$\overbrace{\|\textstyle\sum_i \alpha_i y^i x^i\|^2}$$

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

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

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

Dual Problem : $\displaystyle\min_{\alpha > 0} h(\alpha)$ subject to

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

# 6 Missing Proofs

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

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

↳ Use assumption

$f(w) = \max_{k=1,\ldots,K} f_k(w) \geq f_k(w) \geq \underbrace{f_k(w_0)}_{\Rightarrow f(w_0) \text{ (assumption)}} + \langle v, w-w_0 \rangle = f(w_0) + \langle v, w-w_0 \rangle \quad \square$

$w_t - \eta v$

b) $\|w_{t+1} - w^*\|^2 = \|w_{t+1}\|^2 + \|w^*\|^2 - 2\langle w_{t+1}, w^* \rangle = \|w_t - w^*\|^2$

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

$\|w_t - w^*\|^2 = \|w_t\|^2 + \|w^*\|^2 - 2\langle w_t, w^* \rangle \qquad \leq \|v\| \|w_t - w^*\| \qquad (2)$
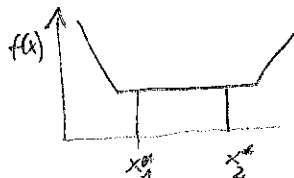
$v \in \partial f(w_t) \Leftrightarrow f(w) \geq f(w_t) + \langle v, w-w_t \rangle \stackrel{w=w^*}{\Longleftarrow} \overbrace{\langle v, w_t - w^* \rangle} \geq f(w_t) - f(w_*) \qquad (3)$

$f$ is convex $\Rightarrow f(\theta w_1 + (1-\theta) w_2) \leq \theta f(w_1) + (1-\theta) f(w_2)$

$\stackrel{w=w_{t+1}}{\Longrightarrow} f(w_t) - f(w_{t+1}) \leq \eta \|v\|^2$

↳ $\|w_{t+1} - w^*\|^2 = \|w_t - \eta v - w^*\| = \|w_t - w^*\|^2 - 2\eta \overbrace{v^T(w_t - w^*)}^{\langle v, w_t - w^* \rangle} + \eta^2 \|v\|^2$

$\stackrel{(3)}{\leq} \|w_t - w^*\|^2 \underbrace{- 2\eta (f(w_t) - f(w^*)) + \eta^2 \|v\|^2}$

$< 0 \quad \text{iff} \quad \eta^2 \|v\|^2 < 2\eta (f(w_t) - f(w^*))$

$\Leftarrow \quad \eta < \frac{2(f(w_t) - f(w^*))}{\|v\|^2} \quad (\geq 0)$

$< \|w_t - w^*\| \quad \text{if} \quad \eta < \frac{2(f(w_t) - f(w^*))}{\|v\|^2} \quad (\geq 0) \quad \square$

c) A convex function can have multiple minima only if it is not strictly convex (i.e. $f(\theta w_1 + (1-\theta) w_2) < \theta f(w_1) + (1-\theta) f(w_2)$ is not guaranteed $\forall w_1, w_2, \theta \in [0,1]$).
If this is the case, multiple global minima can occur, but they are "next to each other", i.e. if $x_1^*$ is a minimum point and $f(x)$ has multiple minima then any other minima $x_2^*$ will be in a ball around $x_1^*$: $x_1^* = \arg\min_K f(x) \wedge x_2^* = \arg\min_K f(x)$

⟹ any convex combination $\theta x_1^* + (1-\theta) x_2^*$ will also be a minimum



d) Since $g(\alpha)$ is the pointwise maximum of a set of convex functions (in fact
$g(\alpha) = \max_\theta h_\theta(\alpha)$ where $h_\theta(\alpha) = f(\theta) + \sum_{i=1}^k \alpha_i g_i(\theta)$ is affine in $\alpha$ ),
$g(\alpha)$ is convex in $\alpha$. For a proof, see Boyd/Vandenberghe Sec. 3.2.3