

収束定理

- **Theorem:** (Novikoff) Let S be a non-trivial training set, and let

$$R = \max_{1 \leq i \leq \ell} \|\mathbf{x}_i\|.$$

Suppose that there exists a vector \mathbf{w}_{opt} such that $\|\mathbf{w}_{opt}\| = 1$ and

$$y_i(\langle \mathbf{w}_{opt}, \mathbf{x}_i \rangle + b_{opt}) \geq \gamma$$

for $1 \leq i \leq \ell$. Then the number of mistakes made by the on-line perceptron algorithm on S is at most

$$\left(\frac{2R}{\gamma}\right)^2.$$

- 間違いの回数に上限 \implies 有限回の繰り返しで収束
- 線形識別可能性を仮定
- 線形識別可能でなければ修正を繰り返し収束しない
- 線形識別可能性の判断には使えない
- 途中で打ち切った場合, そのときの重みが最適かどうかわからない

Define $\hat{\mathbf{x}}_i = (\mathbf{x}_i^T, R)^T$, $\hat{\mathbf{w}} = (\mathbf{w}^T, b/R)^T$. Let $\hat{\mathbf{w}}_{t-1}$ be the augmented weight vector prior to the t th mistake. The t th update is performed when

$$y_i \langle \hat{\mathbf{w}}_{t-1}, \hat{\mathbf{x}}_i \rangle = y_i (\langle \mathbf{w}_{t-1}, \mathbf{x}_i \rangle + b_{t-1}) \leq 0$$

where (\mathbf{x}_i, y_i) is the point incorrectly classified by $\hat{\mathbf{w}}_{t-1}$. The update is the following:

$$\hat{\mathbf{w}}_t = \begin{bmatrix} \mathbf{w}_t \\ b_t/R \end{bmatrix} = \begin{bmatrix} \mathbf{w}_{t-1} \\ b_{t-1}/R \end{bmatrix} + \eta y_i \begin{bmatrix} \mathbf{x}_i \\ R \end{bmatrix} = \hat{\mathbf{w}}_{t-1} + \eta y_i \hat{\mathbf{x}}_i$$

where $b_t = b_{t-1} + \eta y_i R^2$. Since the margin is γ , we have

$$\langle \hat{\mathbf{w}}_t, \hat{\mathbf{w}}_{opt} \rangle = \langle \hat{\mathbf{w}}_{t-1}, \hat{\mathbf{w}}_{opt} \rangle + \eta y_i \langle \hat{\mathbf{x}}_i, \hat{\mathbf{w}}_{opt} \rangle \geq \langle \hat{\mathbf{w}}_{i-1}, \hat{\mathbf{w}}_{opt} \rangle + \eta \gamma$$

and this implies that

$$\langle \hat{\mathbf{w}}_t, \hat{\mathbf{w}}_{opt} \rangle \geq t \eta \gamma$$

Similarly, we have

$$\begin{aligned} \|\hat{\mathbf{w}}_t\|^2 &= \|\hat{\mathbf{w}}_{t-1}\|^2 + 2\eta y_i \langle \hat{\mathbf{w}}_{t-1}, \hat{\mathbf{x}}_i \rangle + \eta^2 \|\hat{\mathbf{x}}_i\|^2 \\ &\geq \|\hat{\mathbf{w}}_{t-1}\|^2 + \eta^2 \|\hat{\mathbf{x}}_i\|^2 = \|\hat{\mathbf{w}}_{t-1}\|^2 + \eta^2 (\|\mathbf{x}_i\|^2 + R^2), \end{aligned}$$

which implies $\|\hat{\mathbf{w}}_t\|^2 \geq \|\hat{\mathbf{w}}_{t-1}\|^2 + \eta^2 R^2$. Thus we have

$$\|\hat{\mathbf{w}}_t\|^2 \geq 2t\eta^2 R^2.$$

The two inequalities combined give the ‘squeezing’ relations

$$\|\mathbf{w}_{opt}\| \sqrt{2t\eta} R \geq \|\mathbf{w}_{opt}\| \|\mathbf{w}_t\| \geq \langle \hat{\mathbf{w}}_t, \hat{\mathbf{w}}_{opt} \rangle \geq t \eta \gamma,$$

which together imply the bound (note $\|\hat{\mathbf{w}}_{opt}\|^2 \leq \|\mathbf{w}_{opt}\|^2 + 1 = 2$)

$$t \leq 2 \left(\frac{R}{\gamma} \right)^2 \|\hat{\mathbf{w}}_{opt}\|^2 \leq \left(\frac{2R}{\gamma} \right)^2.$$

- **Formulation:** (Optimization problem) Given functions f , g_i , $i = 1, \dots, k$, and h_i , $i = 1, \dots, m$, defined on a domain $\Omega \subseteq R^n$, optimization problem is formalized as follows:

$$\begin{aligned} & \text{minimize} && f(\mathbf{w}), && \mathbf{w} \in \Omega \\ & \text{subject to} && g_i(\mathbf{w}) \leq 0, && i = 1, \dots, k \\ & && h_i(\mathbf{w}) = 0, && i = 1, \dots, m \end{aligned}$$

where f is called the *objective function*, and the remaining relations are called, respectively, the *inequality* and *equality constraints*.

- **Definition:** (Convexity) A real-valued function $f(\mathbf{w})$ is called *convex* for $\mathbf{w} \in R^n$ if, $\forall \mathbf{w}, \mathbf{u} \in R^n$, and for any $\theta \in (0, 1)$,

$$f(\theta \mathbf{w} + (1 - \theta) \mathbf{u}) \leq \theta f(\mathbf{w}) + (1 - \theta) f(\mathbf{u})$$

- **Theorem:** (Fermat) A necessary condition for \mathbf{w}^* to be a minimum of a function $f(\mathbf{w})$ is

$$\frac{\partial f(\mathbf{w}^*)}{\partial \mathbf{w}} = \mathbf{0}.$$

This condition, together with convexity of f , is also sufficient.

Lagrange 法

- **Definition:** (Lagrangian) Given an optimization problem with objective function $f(\mathbf{w})$, and equality constraints $h_i(\mathbf{w}) = 0$, $i = 1, \dots, m$, we define the *Lagrangian function* as

$$L(\mathbf{w}, \boldsymbol{\beta}) = f(\mathbf{w}) + \sum_{i=1}^m \beta_i h_i(\mathbf{w})$$

where the coefficients β_i are called the *Lagrange multipliers*.

- **Theorem:** (Lagrange) A necessary condition for a point \mathbf{w}^* to be a minimum of $f(\mathbf{w})$ subject to $h_i(\mathbf{w}) = 0$, $i = 1, \dots, m$, is

$$\begin{aligned} \frac{\partial L(\mathbf{w}^*, \boldsymbol{\beta}^*)}{\partial \mathbf{w}} &= \mathbf{0} \\ \frac{\partial L(\mathbf{w}^*, \boldsymbol{\beta}^*)}{\partial \boldsymbol{\beta}} &= \mathbf{0} \end{aligned}$$

for some values $\boldsymbol{\beta}^*$. The above conditions are also sufficient provided that $L(\mathbf{w}, \boldsymbol{\beta}^*)$ is a convex function of \mathbf{w} .

Kuhn-Tucker 法

- **Definition:** (generalized Lagrangian) Given an optimization problem with domain $\Omega \subseteq R^n$,

$$\begin{aligned} & \text{minimize} && f(\mathbf{w}), && \mathbf{w} \in \Omega \\ & \text{subject to} && g_i(\mathbf{w}) \leq 0, && i = 1, \dots, k \\ & && h_i(\mathbf{w}) = 0, && i = 1, \dots, m \end{aligned}$$

we define the *generalized Lagrangian function* as

$$\begin{aligned} L(\mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta}) &= f(\mathbf{w}) + \sum_{i=1}^k \alpha_i g_i(\mathbf{w}) + \sum_{i=1}^m \beta_i h_i(\mathbf{w}) \\ &= f(\mathbf{w}) + \boldsymbol{\alpha}^T \mathbf{g}(\mathbf{w}) + \boldsymbol{\beta}^T \mathbf{h}(\mathbf{w}) \end{aligned}$$

- **Theorem:** (Kuhn-Tucker) Sufficient conditions for a point \mathbf{w}^* to be an optimum are the existence of $\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*$ such that

$$\begin{aligned} \frac{\partial L(\mathbf{w}^*, \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*)}{\partial \mathbf{w}} &= \mathbf{0}, \\ \frac{\partial L(\mathbf{w}^*, \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*)}{\partial \boldsymbol{\beta}} &= \mathbf{0}, \\ \alpha_i^* g_i(\mathbf{w}^*) &= 0, & i = 1, \dots, k, \\ g_i(\mathbf{w}^*) &\leq 0, & i = 1, \dots, k, \\ \alpha_i^* &\geq 0, & i = 1, \dots, k. \end{aligned}$$

- 例題 1

$$\begin{aligned} & \text{maximize} && x + y \\ & \text{subject to} && x^2 + y^2 \leq 1, \quad x \geq 0, \quad y \geq 0 \end{aligned}$$

- 例題 2

$$\begin{aligned} & \text{maximize} && (x - 1)^2 + (y - 1)^2 \\ & \text{subject to} && x + 2y \leq 1, \quad x \geq 0, \quad y \geq 0 \end{aligned}$$

Kuhn-Tucker 法 (例題)

• 解答 1:

$$\begin{aligned}f(\mathbf{x}) &= -x - y \\g_1(\mathbf{x}) &= x^2 + y^2 - 1 \\g_2(\mathbf{x}) &= -x \\g_3(\mathbf{x}) &= -y\end{aligned}$$

とおく . また $D = \frac{\partial}{\partial \mathbf{x}}$ と書くことにすると ,

$$\begin{aligned}Df(\mathbf{x}) &= (-1, -1) \\Dg_1(\mathbf{x}) &= (2x, 2y) \\Dg_2(\mathbf{x}) &= (-1, 0) \\Dg_3(\mathbf{x}) &= (0, -1)\end{aligned}$$

となる . したがって最適解は

$$\begin{aligned}0 &= (-1, -1) + \alpha_1(2x, 2y) + \alpha_2(-1, 0) + \alpha_3(0, -1) \\0 &= \alpha_1(x^2 + y^2 - 1) \\0 &= \alpha_2(-x) \\0 &= \alpha_3(-y) \\\alpha_i &\geq 0\end{aligned}$$

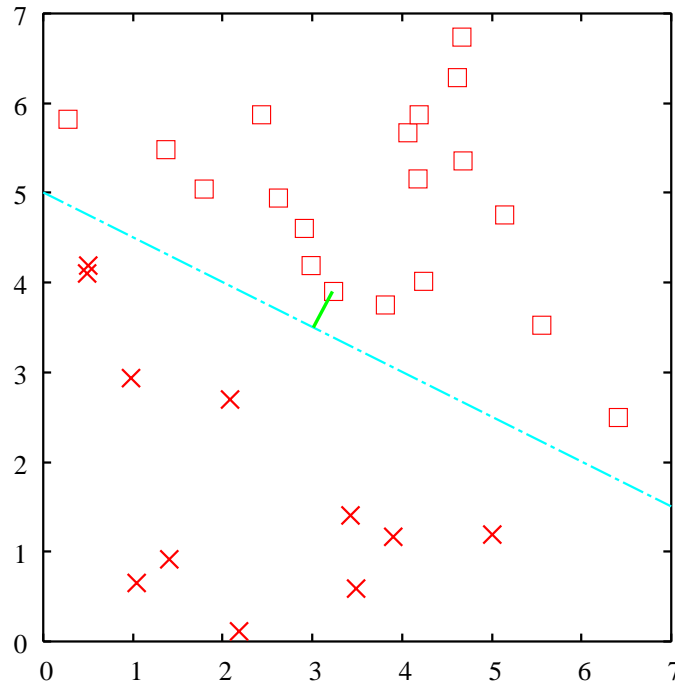
をみます .

1. $x > 0, y > 0$ のとき . $\alpha_2 = \alpha_3 = 0$ であるから $2\alpha_1 x = 2\alpha_1 y = 1$. これを第 2 式に代入し , $\alpha_1 > 0$ を考慮すると $\alpha_1 = 1/\sqrt{2}, x = 1/\sqrt{2}, y = 1/\sqrt{2}$ を得る .
2. $x = 0$ のとき , 第 1 式より $\alpha_2 = -1$ となり不適 .
3. $y = 0$ のとき , 第 1 式より $\alpha_3 = -1$ となり不適 .

したがって , 求める解は $(x, y) = (1/\sqrt{2}, 1/\sqrt{2})$.

- 解答 2: 省略 . 求める解は $(x, y) = (0, 0)$.

- Functional margin: $\gamma = y(\langle \mathbf{w}, \mathbf{x} \rangle + b)$



- Scale \mathbf{w} and b so that: $y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 \quad \forall i$
- Support vectors: $\mathbf{x}^+, \mathbf{x}^-$

$$\langle \mathbf{w}, \mathbf{x}^+ \rangle + b = 1, \quad \langle \mathbf{w}, \mathbf{x}^- \rangle + b = -1$$

- Geometric margin: d

$$\begin{aligned} d &= \frac{1}{2} \left(\left\langle \frac{\mathbf{w}}{\|\mathbf{w}\|}, \mathbf{x}^+ \right\rangle \right) - \left(\left\langle \frac{\mathbf{w}}{\|\mathbf{w}\|}, \mathbf{x}^- \right\rangle \right) \\ &= \frac{1}{2\|\mathbf{w}\|} (\langle \mathbf{w}, \mathbf{x}^+ \rangle) - (\langle \mathbf{w}, \mathbf{x}^- \rangle) = \frac{1}{\|\mathbf{w}\|} \end{aligned}$$

最大マージン識別器 (Primal form)

Proposition: Given a linearly separable training sample

$$S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_\ell, y_\ell))$$

the hyperplane (\mathbf{w}, b) that solves the optimization problem

$$\begin{aligned} & \text{minimize}_{\mathbf{w}, b} \quad \langle \mathbf{w}, \mathbf{w} \rangle \\ & \text{subject to} \quad y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1, \quad \text{for } i = 1, \dots, \ell \end{aligned}$$

realizes the maximal margin hyperplane with geometric margin $\gamma = 1/\|\mathbf{w}\|$.

- Lagrangian

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle - \sum_{i=1}^{\ell} \alpha_i [y_i(\langle \mathbf{w}_i, \mathbf{x}_i \rangle + b) - 1]$$

where $\alpha_i \geq 0$ are Lagrange multipliers. Imposing stationarity condition, we have

$$\begin{aligned} \frac{\partial L(\mathbf{w}, b, \alpha)}{\partial \mathbf{w}} &= \mathbf{w} - \sum_{i=1}^{\ell} y_i \alpha_i \mathbf{x}_i = \mathbf{0}, \\ \frac{\partial L(\mathbf{w}, b, \alpha)}{\partial b} &= \sum_{i=1}^{\ell} y_i \alpha_i = 0. \end{aligned} \tag{1}$$

Substituting these into the primal to obtain

$$L(\mathbf{w}, b, \alpha) = \sum_{j=1}^{\ell} \alpha_j - \frac{1}{2} \sum_{i=1}^{\ell} y_i \alpha_i \sum_{j=1}^{\ell} y_j \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

最大マージン識別器 (Dual form)

Proposition: Given a linearly separable training sample

$$S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_\ell, y_\ell))$$

and suppose the parameters $\boldsymbol{\alpha}^*$ solve the following quadratic optimization problem:

$$\begin{aligned} & \text{maximize} && W(\boldsymbol{\alpha}) = \sum_{j=1}^{\ell} \alpha_j - \frac{1}{2} \sum_{i,j=1}^{\ell} y_i y_j \alpha_i \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\ & \text{subject to} && \sum_{i=1}^{\ell} y_i \alpha_i = 0, \\ & && \alpha_i \geq 0 \text{ for } i = 1, \dots, \ell \end{aligned}$$

Then the weight vector $\mathbf{w} = \sum_{i=1}^{\ell} y_i \alpha_i^* \mathbf{x}_i$ realizes the maximal margin hyperplane with geometric margin $\gamma = 1/\|\mathbf{w}^*\|$.

- **Remark 1:** The value of b

$$b^* = -\frac{1}{2} \left(\max_{y_i=-1} (\langle \mathbf{w}^*, \mathbf{x}_i \rangle) + \min_{y_i=1} (\langle \mathbf{w}^*, \mathbf{x}_i \rangle) \right)$$

- **Remark 2:** Karush-Kuhn-Tucker conditions state that the optimal solutions $\boldsymbol{\alpha}^*$, (\boldsymbol{w}^*, b^*) must satisfy

$$\alpha_i^* [y_i (\langle \boldsymbol{w}^*, \boldsymbol{x}_i \rangle + b^*) - 1] = 0$$

Only for inputs \boldsymbol{x}_i for which the functional margin is one (and therefore lie closest to the hyperplane), the corresponding α_i^* are non-zero. All the other parameters α_i^* are zero.

- **Remark 3:** The optimal hyperplane can be expressed in terms of **support vectors**

$$f(\boldsymbol{x}, \boldsymbol{\alpha}^*, b^*) = \sum_{i=1}^{\ell} y_i \alpha_i \langle \boldsymbol{x}_i, \boldsymbol{x} \rangle + b^* = \sum_{i \in \text{SV}} y_i \alpha_i \langle \boldsymbol{x}_i, \boldsymbol{x} \rangle + b^*$$

Points that are not support vectors have no influence.

- **Remark 4:** Another important consequence of the Karush-Kuhn-Tucker complementarity condition is that for $j \in \text{sv}$,

$$y_j f(\boldsymbol{x}_j, \boldsymbol{\alpha}^*, b^*) = y_j \left(\sum_{i \in \text{SV}} y_i \alpha_i \langle \boldsymbol{x}_i, \boldsymbol{x}_j \rangle + b^* \right) = 1,$$

and therefore

$$\begin{aligned} \langle \boldsymbol{w}^*, \boldsymbol{w}^* \rangle &= \sum_{i,j=1}^{\ell} y_i y_j \alpha_i^* \alpha_j^* \langle \boldsymbol{x}_i, \boldsymbol{x}_j \rangle \\ &= \sum_{j \in \text{SV}} \alpha_j^* y_j \sum_{i \in \text{SV}} y_i \alpha_i^* \langle \boldsymbol{x}_i, \boldsymbol{x}_j \rangle \\ &= \sum_{j \in \text{SV}} \alpha_j^* (1 - y_j b^*) \\ &= \sum_{i \in \text{SV}} \alpha_i^* \end{aligned} \tag{2}$$

i.e.,

$$\gamma = 1/\|\boldsymbol{w}^*\| = \left(\sum_{i \in \text{SV}} \alpha_i^* \right)^{1/2}$$