

Machine Learning

Zhengchao Wan

Peking University

June 13, 2016

Overview

- Geometry of support vector machine
- Boosting

Linear classifier

Given input pattern $\mathbf{x} \in R^n$, consider a linear function

$$f(\mathbf{x}, \mathbf{w}) = \mathbf{w} \cdot \mathbf{x} + b \quad (1)$$

having parameters $\xi = (\mathbf{w}, b)$. The machine classifies patterns into two classes C_+ and C_- , according to the signature of output function f . That is, when $f(\mathbf{x}, \xi) > 0$, \mathbf{x} is classified into C_+ , and otherwise into C_- .

Linearly separable

Training set $D = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ are divided into two classes. When $\mathbf{x}_i \in C_+$, it's given a teacher signal $y_i = 1$, and when $\mathbf{x}_i \in C_-$, $y_i = -1$. They are **linearly separable** when there exists \mathbf{w} and b , for which

$$\mathbf{w} \cdot \mathbf{x} + b > 0, \mathbf{x} \in C_+, \mathbf{w} \cdot \mathbf{x} + b < 0, \mathbf{x} \in C_- \quad (2)$$

holds. When (\mathbf{w}, b) is such a solution, also is $(c\mathbf{w}, cb), \forall c > 0$. We can therefore impose the constraints

$$|\mathbf{w} \cdot \mathbf{x} + b| \geq 1, \min_i |\mathbf{w} \cdot \mathbf{x} + b| = 1. \quad (3)$$

Support vector

The Euclidean distance from point \mathbf{x} to the separating hyperplane $\mathbf{w} \cdot \mathbf{x} + b = 0$ is

$$d = \frac{|\mathbf{w} \cdot \mathbf{x} + b|}{|\mathbf{w}|} \quad (4)$$

The minimum of distances $d_i = \frac{|\mathbf{w} \cdot \mathbf{x}_i + b|}{|\mathbf{w}|}$ is given by $d_{min} = \frac{1}{|\mathbf{w}|}$ and is attained by the points \mathbf{x}_i that satisfy

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1. \quad (5)$$

We call these points the **support vectors** of the training set D and the minimal distance the **margin**.

Support vector

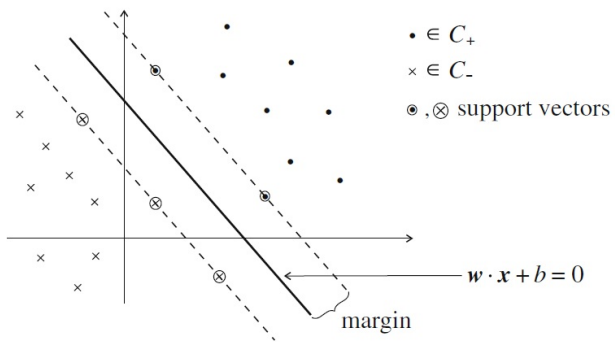


Figure: Linear classifier and support vectors

Optimal linear machine

A good machine has a large margin so we want to minimize

$$C(\mathbf{w}) = \frac{1}{2}|\mathbf{w}|^2 \quad (6)$$

under the constraint

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1. \quad (7)$$

Using Lagrange multipliers $\alpha = (\alpha_1, \dots, \alpha_N)$, we transform the problem into the unconstrained minimization of

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2}|\mathbf{w}|^2 - \sum \alpha_i (y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1). \quad (8)$$

Optimal linear machine

By differentiating with \mathbf{w} and b and making the derivatives equal to 0, we have $\sum_i \alpha_i y_i = 0$, $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$ and the dual problem

$$\begin{aligned} \text{maximize } L^*(\boldsymbol{\alpha}) &= \sum \alpha_i - \frac{1}{2} \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j, \\ \alpha_i &\geq 0, \sum \alpha_i y_i = 0. \end{aligned} \quad (9)$$

This is a quadratic function of α_i , which has a specific algorithm to solve. It should be remarked that $\alpha_i = 0$ when \mathbf{x}_i is not a support vector.

The optimized output function is written as

$$f(\mathbf{x}, \mathbf{w}) = \sum \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b \quad (10)$$

in terms of the solution α_i , which only uses the support vectors.

Embedding into high-dimensional space

Consider a nonlinear transformation of $\mathbf{x} \in R^n$ into a high-dimensional space $R^m (m > n)$ by

$$\mathbf{z} = \varphi(\mathbf{x}) = [\varphi_1(\mathbf{x}), \dots, \varphi_m(\mathbf{x})]. \quad (11)$$

The original classification problem is transformed into a linear classification problem in R^m , where

$$f(\mathbf{x}, \boldsymbol{\xi}) = \mathbf{w} \cdot \varphi(\mathbf{x}) + b, \quad \boldsymbol{\xi} = (\mathbf{w}, b). \quad (12)$$

Embedding into high-dimensional space

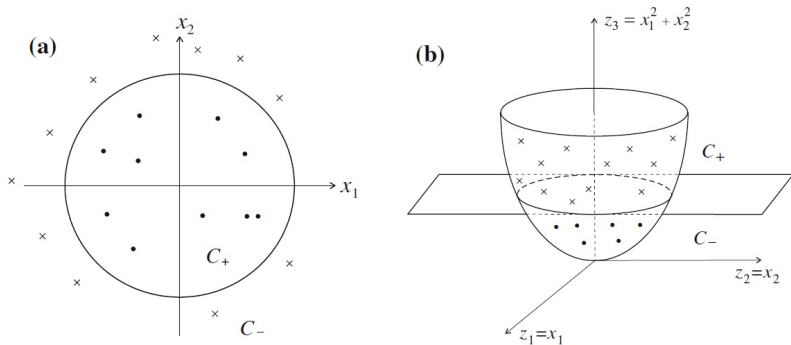


Fig. 11.7 a Non-separable in \mathbb{R}^2 ; b separable in \mathbb{R}^3

Kernel method

Consider the inner product of $\mathbf{z} = \varphi(\mathbf{x})$ and $\mathbf{z}' = \varphi(\mathbf{x}')$ after embedding,

$$K(\mathbf{x}, \mathbf{x}') = \mathbf{z}(\mathbf{x}) \cdot \mathbf{z}(\mathbf{x}') = \sum z_i(\mathbf{x})z_i(\mathbf{x}') \quad (13)$$

We call K as a kernel and then the optimal output function can be written as

$$f(\mathbf{x}, \mathbf{w}) = \sum \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b, \quad (14)$$

which means we don't need to compute the inner product of φ .

Kernel method

We may start from a kernel function without specifying embedding functions, provides $K(\mathbf{x}, \mathbf{x}')$ is positive-definite satisfying

$$\sum c_i c_j K(\mathbf{x}_i, \mathbf{x}_j) > 0, \forall \mathbf{c} = (c_1, \dots, c_l), \quad (15)$$

called the Mercer condition.

- **Gaussian kernel**

$$K(\mathbf{x}, \mathbf{x}') = \exp\left\{-\frac{|\mathbf{x} - \mathbf{x}'|^2}{\sigma^2}\right\} \quad (16)$$

- **Polynomial kernel**

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}' + 1)^p \quad (17)$$

Riemannian metric induced by
kernel

The original space R^n of patterns is embedded in R^m , possibly in R^∞ , as a curved n -dimensional submanifold.

$$ds^2 = |\varphi(\mathbf{x} + d\mathbf{x}) - \varphi(\mathbf{x})|^2 = \sum \frac{\partial}{\partial x_i} \varphi(\mathbf{x}) \cdot \frac{\partial}{\partial x_j} \varphi(\mathbf{x}) dx_i dx_j \quad (18)$$

Therefore, the induced Riemannian metric is given by

$$g_{ij}(\mathbf{x}) = \left(\frac{\partial}{\partial x_i} \varphi(\mathbf{x}) \right) \cdot \left(\frac{\partial}{\partial x_j} \varphi(\mathbf{x}) \right) \quad (19)$$

Riemannian metric induced by
kernel

Rewrite it in terms of the kernel as

$$g_{ij}(\mathbf{x}) = \frac{\partial^2}{\partial x_i \partial x'_j} K(\mathbf{x}, \mathbf{x}') \Big|_{\mathbf{x}' = \mathbf{x}} \quad (20)$$

The volume element at point \mathbf{x} is given by

$$dV(\mathbf{x}) = \sqrt{|g_{ij}(\mathbf{x})|} dx_1 \cdots dx_n, \quad (21)$$

which shows how the volume is enlarged or contracted at around \mathbf{x} .

Geometric insight

To increase the margin or separability of classes, we need to enlarge the spatial resolution around the boundary $f(\mathbf{x}) = 0$, which is expressed by

$$ds^2 = \sum g_{ij} dx_i dx_j.$$

In practice, we don't know the boundary, so by using the knowledge that support vectors are (mostly) located around the boundary, we solve the problem by increasing the metric in the neighborhood of the support vectors.

Conformal transformation

A conformal transformation

$$\tilde{g}_{ij}(\mathbf{x}) = \sigma(\mathbf{x})g_{ij}(\mathbf{x}), \quad (22)$$

gives a solution to this problem, for the metric is enlarged by a factor $\sigma(\mathbf{x})$, which should have large values at the SV positions. In practice, it's difficult to find a mapping $\tilde{\varphi}$ which brings the above conformal transformation. Therefore, we consider a quasi-conformal transformation obtained from modification of a kernel.

Conformal transformation of kernel

Definition (Conformal transformation of kernel)

Let $\sigma(\mathbf{x})$ be a positive scalar function.

$$\tilde{K}(\mathbf{x}, \mathbf{x}') = \sigma(\mathbf{x})\sigma(\mathbf{x}')K(\mathbf{x}, \mathbf{x}') \quad (23)$$

is called a conformal transformation of a kernel by factor $\sigma(\mathbf{x})$.

In Amari and Wu(1999),

$$\sigma(\mathbf{x}) = \sum_i e^{-\kappa_i |\mathbf{x} - \mathbf{x}_i^*|} \quad (24)$$

was chosen, where \mathbf{x}_i^* are the SV and κ_i are adequate constants.

Conformal transformation of kernel

$$\begin{aligned}\tilde{g}_{ij}(\mathbf{x}) &= \sigma^2(\mathbf{x})g_{ij}(\mathbf{x}) + \sigma_i(\mathbf{x})\sigma_j(\mathbf{x})K(\mathbf{x}, \mathbf{x}) \\ &\quad + \sigma(\mathbf{x})\{\sigma_i(\mathbf{x})K_j(\mathbf{x}, \mathbf{x}) + \sigma_j(\mathbf{x})K_i(\mathbf{x}, \mathbf{x})\},\end{aligned}\quad (25)$$

where

$$\sigma_i = \frac{\partial}{\partial x_i}\sigma(\mathbf{x}), \quad K_i(\mathbf{x}, \mathbf{x}) = \frac{\partial}{\partial x_i}K(\mathbf{x}, \mathbf{x}')|_{\mathbf{x}'=\mathbf{x}}. \quad (26)$$

When $K_i(\mathbf{x}, \mathbf{x}) = 0$, which is satisfied by the Gaussian kernel, we have a simplified expression

$$\tilde{g}_{ij}(\mathbf{x}) = \sigma^2(\mathbf{x})g_{ij}(\mathbf{x}) + \sigma_i(\mathbf{x})\sigma_j(\mathbf{x})K(\mathbf{x}, \mathbf{x}). \quad (27)$$

Boosting

Boosting refers to a general and provably effective method of producing a very accurate prediction rule by combining rough and moderately inaccurate rules of thumb.

Boosting

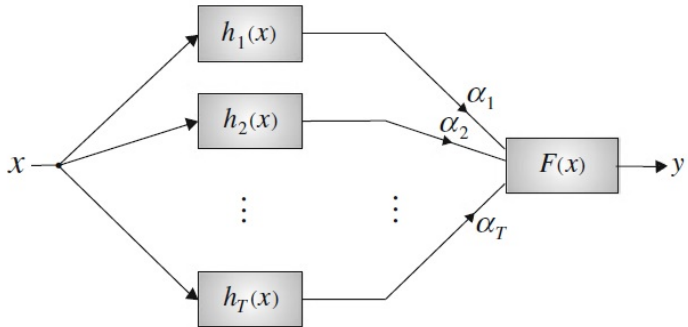


Figure: Integration of weak machines

AdaBoost

The **AdaBoost** algorithm is a famous algorithm in solving the boosting problem. The algorithm takes as input a training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$ where each \mathbf{x}_i belongs to some domain \mathcal{X} and each **label** y_i is in some label set \mathcal{Y} , and mostly, we use $\mathcal{Y} = \{-1, +1\}$.

AdaBoost generates weak machines in a series of rounds $t = 1, \dots, T$. One main idea of the algorithm is to maintain a distribution or set of weights over the training set.

AdaBoost

Algorithm (AdaBoost)

Given $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$, where $\mathbf{x}_i \in \mathcal{X}, y_i \in \mathcal{Y} = \{-1, +1\}$ Initialize $D_1(i) = 1/m$. For $t = 1, \dots, T$:

- Train weak learner h_t using distribution D_y .
- Choose $\alpha_t = \frac{1}{2} \ln\left(\frac{1-\epsilon_t}{\epsilon_t}\right)$, where $\epsilon_t = P_{i \sim D_t}[h_t(\mathbf{x}_i) \neq y_i]$.
- Update:

$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} \times \begin{cases} e^{-\alpha_t} & \text{if } h_t(\mathbf{x}_i) = y_i \\ e^{\alpha_t} & \text{if } h_t(\mathbf{x}_i) \neq y_i \end{cases}$$

$$= \frac{D_t(i) \exp(-\alpha_t y_i h_t(\mathbf{x}_i))}{Z_t}$$

where Z_t is a normalization factor.

Output the final hypothesis: $H(\mathbf{x}) = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(\mathbf{x})\right)$.

Measure manifold and probability
manifold

$$\mathcal{M} = \left\{ m(y|\mathbf{x}) \mid \sum_{y \in \mathcal{Y}} m(y|\mathbf{x}) < \infty (a.e. \mathbf{x}) \right\} \quad (28)$$

$$\mathcal{P} = \left\{ m(y|\mathbf{x}) \mid \sum_{y \in \mathcal{Y}} m(y|\mathbf{x}) = 1 (a.e. \mathbf{x}) \right\} \quad (29)$$

Empirical conditional probability

For given examples $\{(\mathbf{x}_i, y_i); i = 1, \dots, n\}$, let

$$\tilde{p}(y|\mathbf{x}) = \begin{cases} \delta(y_i, y), & \mathbf{x} = \mathbf{x}_i, \\ \frac{1}{|\mathcal{Y}|}, & \text{otherwise} \end{cases} \quad (30)$$

be the empirical conditional probability density. Suppose a unique label y_i is given for each input \mathbf{x}_i .

KL divergence

For two points $p, q \in \mathcal{M}$, the KL divergence extended over \mathcal{M} is defined by

$$D(p, q) = \int_{\mathcal{X}} \mu(\mathbf{x}) \sum_{y \in \mathcal{Y}} \left(p(y|\mathbf{x}) \log \frac{p(y|\mathbf{x})}{q(y|\mathbf{x})} - p(y|\mathbf{x}) + q(y|\mathbf{x}) \right) d\mathbf{x} \quad (31)$$

where $\mu(\mathbf{x})$ is the marginal distribution of \mathbf{x} , and in the most cases of the following discussion, we fix $\mu(\mathbf{x})$ with the empirical distribution $\mu(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \delta(\mathbf{x}_i, \mathbf{x})$.

Probability families

Given a set functions $\mathbf{f} = \{f_t(\mathbf{x}, y); t = 1, \dots, T\}$ on $\mathcal{X} \times \mathcal{Y}$ and let $\tilde{f}(\mathbf{x}) = E_{\tilde{p}}(f|\mathbf{x}) = \sum_{y \in \mathcal{Y}} \tilde{p}(y|\mathbf{x})f(\mathbf{x}, y)$.

e-flat subspace

$$\begin{aligned} \mathcal{Q} &= \mathcal{Q}(q_0, \tilde{p}, \mathbf{f}) \\ &= \left\{ q \in \mathcal{M} \mid q(y|\mathbf{x}) = q_0(y|\mathbf{x}) \exp\left(\sum_{t=1}^T \alpha_t (f_t(\mathbf{x}, y) - \tilde{f}_t(\mathbf{x})) \right) \right\} \end{aligned}$$

m-flat subspace

$$\begin{aligned} \mathcal{F} &= \mathcal{F}(\tilde{p}, \mathbf{f}) \\ &= \left\{ q \in \mathcal{M} \mid \int_{\mathcal{X}} \mu(\mathbf{x}) \sum_{y \in \mathcal{Y}} q(y|\mathbf{x}) (f_t(\mathbf{x}, y) - \tilde{f}_t(\mathbf{x})) d\mathbf{x} = 0; \forall t \right\} \end{aligned}$$

Flatness

$$\tilde{p}(y|\mathbf{x}) \in \mathcal{F}.$$

We can conclude the flatness of the above two families from the following equations.

$$q_0 \exp\left(\beta_1 \log \frac{q_1}{q_0} + \beta_2 \log \frac{q_2}{q_0}\right) \in \mathcal{Q}$$

$$\beta_1 p_1 + \beta_2 p_2 \in \mathcal{F}$$

$$\forall \beta_1, \beta_2 > 0, p_1, p_2 \in \mathcal{F}, q_1, q_2 \in \mathcal{Q}$$

Geometric structure of KL
divergence

We know $D = D^{(-1)}$, therefore we have

$$g_{ij}(\xi) = \int \partial_i l^{(-1)}(x; \xi) \partial_j l^{(1)}(x; \xi) dx$$

where $l^{(-1)}(x; \xi) = p(x, \xi)$, $l^{(1)}(x; \xi) = \log p(x; \xi)$.

Hence for any $X, Y \in T_\xi(\mathcal{M})$, $\langle X, Y \rangle_\xi = \int X l^{(-1)} Y l^{(1)} dx$.

Orthogonality

Suppose $q^* \in \mathcal{F} \cap \mathcal{Q}$. For any $p \in \mathcal{F}, q \in \mathcal{Q}$, let

$$\gamma_1^{(-1)}(x; t) = tq^{(-1)} + (1-t)q^{*(-1)} \text{ and}$$

$\gamma_2^{(1)}(x; t) = tp^{(1)} + (1-t)q^{*(1)}$ be the m-geodesic connecting p, q^* and the e-geodesic connecting q, q^* respectively.

$$\begin{aligned} \langle \dot{\gamma}_1(0), \dot{\gamma}_2(0) \rangle_{q^*} &= \int \frac{d\gamma_1^{(-1)}(x; t)}{dt} \Big|_{t=0} \frac{d\gamma_2^{(1)}(x; s)}{ds} \Big|_{s=0} dx \\ &= \int (q^{(-1)} - q^{*(-1)})(p^{(1)} - q^{*(1)}) dx \\ &= \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} \left(\sum_{t=1}^T \beta_t (f_t(\mathbf{x}, y) - \tilde{f}_t(\mathbf{x})) \right) (q(y|\mathbf{x}) - q^*(y|\mathbf{x})) d\mu(\mathbf{x}) = 0 \end{aligned}$$

Thus, \mathcal{F} and \mathcal{Q} is orthogonal at q^* .

Pythagorean relation

Lemma

For any $p \in \mathcal{F}(\tilde{p}, \mathbf{f})$ and $q \in \mathcal{Q}(q_0, \tilde{p}, \mathbf{f})$, the Pythagorean relation

$$D(p, q) = D(p, q^*) + D(q^*, q) \quad (32)$$

holds.

We have immediately that $\{q^*\} = \mathcal{F} \cap \mathcal{Q}$.

Theorem (Lebanon and Lafferty (2001))

$$q^* = \underset{p \in \mathcal{F}(\tilde{p}, \mathbf{f})}{\operatorname{argmin}} D(p, q_0) = \underset{q \in \mathcal{Q}(q_0, \tilde{p}, \mathbf{f})}{\operatorname{argmin}} D(\tilde{p}, q) \quad (33)$$

Pythagorean relation

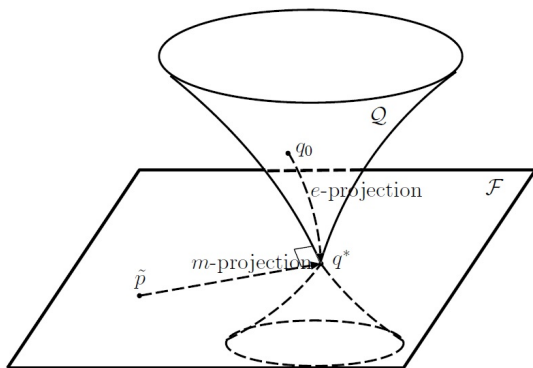


Figure: A geometrical interpretation of the dual optimization problems in the AdaBoost algorithm.

Geometric view of AdaBoost

$$\begin{aligned} Q_t &= Q(q_{t-1}, \tilde{p}_t, f_t) \\ &= \left\{ q \in \mathcal{M} \mid q(y|\mathbf{x}) = q_{t-1}(y|\mathbf{x}) \exp(\alpha_t (f_t(\mathbf{x}, y) - \tilde{f}_t(\mathbf{x}))) \right\} \end{aligned}$$

$$\begin{aligned} \mathcal{F}_t &= \mathcal{F}(\tilde{p}, f_t) \\ &= \left\{ p \in \mathcal{M} \mid \int_{\mathcal{X}} \mu(\mathbf{x}) \sum_{y \in \mathcal{Y}} p(y|\mathbf{x}) (f_t(\mathbf{x}, y) - \tilde{f}_t(\mathbf{x})) d\mathbf{x} = 0 \right\} \end{aligned}$$

It can be verified that Q_t and \mathcal{F}_t intersect at one point q_t .

Geometric view of AdaBoost

Algorithm (AdaBoost)

- ① Initialize $q_0(y|\mathbf{x}) = 1$.
- ② For $t = 1, \dots, T$
 - Select a machine h_t s.t.

$$\sum_{i=1}^n \sum_{y \in \mathcal{Y}} q_{t-1}(y|\mathbf{x}_i) (f_t(\mathbf{x}_i, y) - f_t(\mathbf{x}_i, y_i)) \neq 0$$

$$\text{where } f_t(\mathbf{x}, y) = \begin{cases} 1/2, & y \in h_t(\mathbf{x}), \\ -1/2, & \text{otherwise.} \end{cases}$$

- Construct Q_t and \mathcal{F}_t with \tilde{p}, q_{t-1}, f_t .
- Find q_t and corresponding α_t which is the intersection of Q_t and \mathcal{F}_t

Geometric view of AdaBoost

Algorithm (AdaBoost)

- ③ Output the final decision as the majority vote of $\{h_t; t = 1, \dots, T\}$ with weights $\{\alpha_t; t = 1, \dots, T\}$

$$H(\mathbf{x}) = \operatorname{argmax}_{y \in \mathcal{Y}} \sum_{t=1}^T \alpha_t f_t(\mathbf{x}, y)$$

For the binary case, $\alpha_t = \frac{1}{2} \log \frac{1-\epsilon_t}{\epsilon_t}$, where ϵ_t is the weighted error defined by $\epsilon_t = \sum_{i=1}^n I(y_i \neq h_t(\mathbf{x}_i)) D_t(i)$, $D_t(i) = \frac{q_{t-1}(y \neq y_i | \mathbf{x}_i)}{Z_t}$, and Z_t is a normalization constant.

Geometric view of AdaBoost

In fact since the relation $D(\tilde{p}, q_{t-1}) = D(\tilde{p}, q_t) + D(q_t, q_{t-1})$ holds, as t increases q_t approaches to \tilde{p} as long as $D(q_t, q_{t-1}) > 0$.

$$\begin{aligned} H(\mathbf{x}) &= \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \sum_{t=1}^T \alpha_t f_t(\mathbf{x}, y) \\ &= \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \log q_T(y|\mathbf{x}) \\ &= \underset{y \in \mathcal{Y}}{\operatorname{argmax}} q_T(y|\mathbf{x}) \end{aligned}$$

Geometric view of AdaBoost

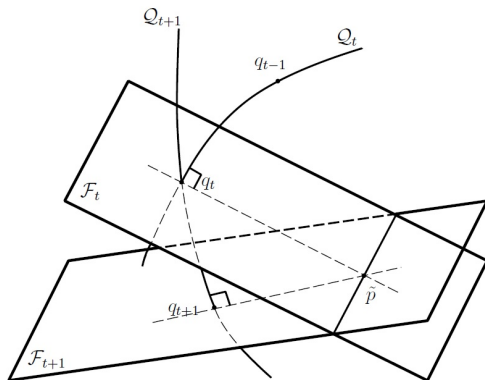


Figure: A geometrical interpretation of the sequential update of AdaBoost.

U -boost

Using Bregman divergence of convex function U (in the KL divergence, $U(z) = \exp(z)$) $D_U(p, q)$ instead of KL divergence, we can obtain a series of new boost algorithms.