

-5

$$L(w, b, \epsilon, \Lambda) = \frac{1}{2} w^T w + \frac{c}{2} \sum_{i=1}^l \epsilon_i^2 - \sum_{i=1}^l \alpha_i \{y_i (w^T x_i + b) - 1 + \epsilon_i\} - \sum_{i=1}^l r_i \epsilon_i$$

$$\frac{\partial L}{\partial w} = 0 \rightarrow \left(w - \sum_{i=1}^l \alpha_i y_i x_i = 0 \right)^{(*)}, \quad \frac{\partial L}{\partial b} = 0 \rightarrow \left(\sum_{i=1}^l \alpha_i y_i = 0 \right)^{(*)}$$

$$\frac{\partial L}{\partial \epsilon_i} = 0 \rightarrow c \epsilon_i - \alpha_i - r_i = 0 \rightarrow c \epsilon_i = r_i + \alpha_i^{(\heartsuit)}$$

$$r_i \epsilon_i = 0^{(+)}, \quad \alpha_i \{y_i (w^T x_i + b) - 1 + \epsilon_i\} = 0$$

$$\Rightarrow L(w, b, \epsilon, \Lambda) = \frac{1}{2} w^T w + \frac{c}{2} \sum_{i=1}^l \epsilon_i^2 - \underbrace{\sum_{i=1}^l \alpha_i y_i w^T x_i}_{-w^T w} - \underbrace{\sum_{i=1}^l \alpha_i y_i b}_0 + \sum_{i=1}^l \alpha_i$$

$$-\sum_{i=1}^l \alpha_i \epsilon_i - \sum_{i=1}^l r_i \epsilon_i =$$

$$-\sum_{i=1}^l (\alpha_i + r_i) \epsilon_i$$

$c \epsilon_i$ طبق (\heartsuit)

$$-\sum_{i=1}^l c \epsilon_i^2$$

$$= \frac{1}{2} w^T w + \frac{c}{2} \sum_{i=1}^l \epsilon_i^2 - w^T w + \sum_{i=1}^l \alpha_i - \sum_{i=1}^l c \epsilon_i^2$$

$$= \left\{ -\frac{1}{2} w^T w - \frac{c}{2} \sum_{i=1}^l \epsilon_i^2 + \sum_{i=1}^l \alpha_i \right\}^{(\bullet)}$$

$c \sum_{i=1}^l \epsilon_i$

$$\begin{aligned} (\heartsuit) &\rightarrow c \epsilon_i = r_i + \alpha_i \rightarrow r_i = c \epsilon_i - \alpha_i \\ (+) &\rightarrow r_i \epsilon_i = 0 \div \epsilon_i \rightarrow r_i = 0 \end{aligned} \quad \left. \vphantom{\begin{aligned} (\heartsuit) &\rightarrow c \epsilon_i = r_i + \alpha_i \\ (+) &\rightarrow r_i \epsilon_i = 0 \end{aligned}} \right\} c \epsilon_i - \alpha_i = 0 \Rightarrow \alpha_i = c \epsilon_i$$

$$(\bullet) \rightarrow L(w, b, \epsilon, \Lambda) = -\frac{1}{2} w^T w - \frac{c}{2} \sum_{i=1}^l \epsilon_i^2 + c \sum_{i=1}^l \epsilon_i$$

6- الف) در روش های **discriminative**، مدل از احتمال شرطی $P(w_i | x)$ برای پیش بینی داده های **unseen** استفاده می کند. این روش ها برای بدست آوردن احتمال، به طور مستقیم یک توزیع فرضی برای $P(w_i | x)$ در نظر نمی گیرند و سپس پارامترهای این توزیع را با کمک داده های **train** تخمین می زنند. مدل های **discriminative** در واقع مرز تصمیم بین کلاس های مختلف را مدل می کنند. برخی مدل های **discriminative** عبارتند از: **logistic Regression**، **SVMs**، **nearest neighbor**، **decision tree & Random Forest** و ... این مدل ها قابلیت تولید داده ی جدید ندارند و به طور کلی هدف این مدل ها، جداسازی یک کلاس داده ها از دیگری است.

روش های **Generative**، برای بدست آوردن احتمال شرطی $P(Y|X)$ ، احتمال $P(Y)$ prior و **likelihood** $P(X|Y)$ را به کمک داده ی **train** تخمین می زنند و سپس طبقه ایابی می کنند.

احتمال **joint** (یعنی $P(x, y)$) را بدست می آورند و این احتمال **joint** می تواند به عنوان جایگزین برای احتمال شرطی $P(Y|X)$ استفاده شود تا به کمک آن پیش بینی انجام دهیم.

$$P(x, y) = P(y) \cdot P(x|y)$$

برخی مدل های **generative** عبارتند از: **Bayesian networks**، **Hidden Markov model**، **GANs** و ... این مدل ها قابلیت تولید داده ی جدید دارند و به طور کلی هدف آن ها مدل سازی خود داده ها است (برخلاف مدل های **discriminative** که هدف آن ها مدل سازی مرز تصمیم بین کلاس های مختلف بود).

در کل، از آنجایی که (همانطوری که توضیح دادم) بدست آوردن $P(x, y)$ بسیار زمانبر است، روش های **discriminative** در مقابل روش های **generative**، **efficient** تر هستند اما نسبت به روش های **generative**، دقت کمتری هم دارند.

ب) مهمترین مزیت روش **one vs. all** این است که مسأله ی طبقه بندی داده های c کلاس را به c تا مسأله ی **binary classification** تبدیل می کند (تبدیل مسأله ی **multi-class classification** به تعدادی مسأله ی ساده تر **binary classification**). اما عیب این روش این است که مجبوریم برای هر کلاس از داده ها، یک مدل بسازیم. ساختن تعداد زیادی مدل مخصوصاً برای زمانی که دیتاست بزرگی داریم و زمانی که تعداد کلاس های داده ی ما زیاد است، با استفاده از مدل های کندی مثل شبکه های عصبی، بسیار چالش برانگیز خواهد بود. نکته ی دیگر آنست که ممکن است توزیع داده ها، کلاس در هر یک از این مسائل **binary classification**، نامتوازن شود و در چنین حالتی طبقه بندی کننده عملکرد خوبی برای کلاس هایی که داده های آن از بقیه کمتر و آشفته تر است نداشته باشد.

ج) هم‌تکین‌ترین dual نسبت به primal زمانی است که می‌خواهیم از kernel برای classify کردن داده‌هایی استفاده کنیم که در فضای اولیه به صورت خطی قابل جداسازی نیستند. dual form به راحتی با استفاده از kernel قابل اعمال است. در چنین حالتی (وقتی از kernel استفاده می‌کنیم)، مسئله‌ی optimization توسط فرم dual به صورت زیر می‌باشد:

$$\max \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \alpha_i \alpha_j k(x_i, x_j)$$

$$\text{s.t. } \forall i \alpha_i \geq 0 \wedge \sum_{i=1}^n y_i \alpha_i = 0$$

می‌بینیم که این فرم تقریباً با فرم اولیه‌ی dual یکسان است. تنها تفاوت این است که به جای dot product عادی، باید تابع kernel را جای‌سبب کنیم. دقت کنید که چنین چیزی در فرم primal ممکن نیست. چرا که باید صریحاً mapping هر نقطه را جای‌سبب کنیم:

$$\min_{w, w_0} \frac{1}{2} w^T w$$

$$\text{s.t. } \forall i: y_i (w^T \phi(x_i) + w_0) \geq 1$$

همچنین بارونش dual، و classify کردن یک نقطه‌ی جدید ساده‌تر خواهد بود.

در کل با استفاده از روش dual، داده‌هایی که به صورت خطی قابل جداسازی نیستند را به راحتی می‌توان classify کرد (با کمی پیچیدگی محاسباتی بیشتر). (اما برای linear SVM، حل کردن فرم primal ساده‌تر است)

-7

$$K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2}\right) = \varphi^t(x_i) \varphi(x_j)$$

(الف)

$$d(\varphi(x_i), \varphi(x_j)) = \|\varphi(x_i) - \varphi(x_j)\| = \sqrt{(\varphi(x_i) - \varphi(x_j))^t (\varphi(x_i) - \varphi(x_j))}$$

$$= \sqrt{\underbrace{\varphi^t(x_i) \varphi(x_i)}_{K(x_i, x_i)} - \underbrace{\varphi^t(x_i) \varphi(x_j)}_{K(x_i, x_j)} - \underbrace{\varphi^t(x_j) \varphi(x_i)}_{= \varphi(x_i)^t \varphi(x_j) = K(x_i, x_j)} + \underbrace{\varphi^t(x_j) \varphi(x_j)}_{K(x_j, x_j)}}$$

$$= \sqrt{K(x_i, x_i) - 2K(x_i, x_j) + K(x_j, x_j)} =$$

 $e^i = 1$ $e^o = 1$

$$= \sqrt{2 - 2K(x_i, x_j)} = \sqrt{2 - 2\exp\left(-\frac{\|x_i - x_j\|^2}{2}\right)} \leq \sqrt{2}$$

$$\Rightarrow \|\varphi(x_i) - \varphi(x_j)\|^2 \leq (\sqrt{2})^2 = 2$$

(ب)

$$1) K(x, y) = f(x) k_1(x, y) f(y)$$

طبق خواص kernel های دایم: $K(x, y) = f(x) \cdot f(y)$ برای هر تابع f روی X ، زیرا فقط

یک ویژگی توسط f define می شود. پس می توانیم بنویسیم:

$$K_2(x, y) = f(x) f(y)$$

(*)

همچنین طبق خواص kernel های دایم: $K(x, y) = k_1(x, y) k_2(x, y)$

$$K(x, y) = k_1(x, y) k_2(x, y) \stackrel{(*)}{=} k_1(x, y) f(x) f(y)$$

پس می توانیم بنویسیم: