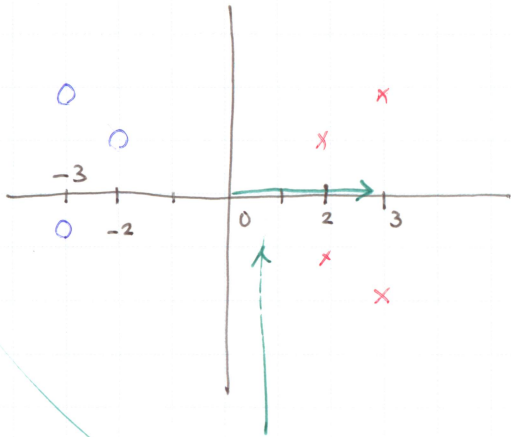


# SVM: Mathematics Behind Large Margin Classification

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

s.t.  $\|\theta\| \cdot p^{(i)} \geq 1$  if  $y^{(i)} = 1$   
 such that  $\|\theta\| \cdot p^{(i)} \leq -1$  if  $y^{(i)} = 0$

$p^{(i)}$  は (符号を正の)  $x^{(i)}$  から  $\theta$  への射影 (projection)



$$\theta = (\theta_1, 0)$$

$$2 \cdot \theta_1 + 1 \cdot 0 \geq 1 \quad \therefore \theta_1 \geq \frac{1}{2}$$

$$-2 \cdot \theta_1 + 1 \cdot 0 \leq -1 \quad \therefore \theta_1 \geq \frac{1}{2}$$

$\|\theta\|$  を最小化したいから

$$\theta = \left(\frac{1}{2}, 0\right)$$

$$\|\theta\| = \frac{1}{2}$$

この方向の vector の長さを最小に求めたい。  
(条件を満たすこと)

# SVM : Kernels

## Non-linear Decision Boundary

$$\text{predict } z=1 \text{ if } \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 + \dots \geq 0$$

kernel ... Given  $x$ , landmark  $l^{(1)}, l^{(2)}, l^{(3)}$  2次元以上の新しく feature を計算せよ.

$$f_1 = \underbrace{\text{similarity}(x, l^{(1)})}_{\uparrow \text{kernel}} = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

$\sum_{j=1}^n (x_j - l_j^{(1)})^2$

$$f_2 = \quad \quad \quad (2) \quad \quad \quad (2) \quad \quad \quad \uparrow \text{Gaussian kernel}$$
$$f_3 = \quad \quad \quad (3) \quad \quad \quad (3)$$

if  $x \approx l^{(1)}$   $f_1 \approx \exp\left(-\frac{0^2}{2\sigma^2}\right) \approx 1$

if  $x$  は  $l^{(1)}$  から遠ければ

$$f_1 = \exp\left(-\frac{\text{大定数}^2}{2\sigma^2}\right) \approx 0$$

# SVM: Kernels I の Quiz

$x_1$  は1次元の feature

$$l^{(1)} = 5$$

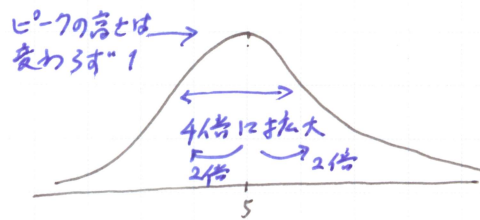
$$f_1 = \exp\left(-\frac{\|x_1 - l^{(1)}\|^2}{2\sigma^2}\right)$$

$\sigma^2 = 1$  の時のグラフ



$$y = e^{-\frac{x}{2}}$$

$\sigma^2 = 4$  のときのグラフ



$$y = e^{-\frac{x}{8}}$$

## SVM: kernels II

Given  $x$ ,  $f_i = \text{similarity}(x, l^{(i)}) = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right)$

Predict  $y=1$ , if  $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$

$l^{(1)}, l^{(2)}, l^{(3)}, \dots$  をどうやって見つければよいか?

### SVM with kernels

$(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$  が与えられた時  
 $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(m)} = x^{(m)}$  とする。

example  $x$  が与えられたと

$$\begin{matrix} f_1 = \text{similarity}(x, l^{(1)}) \\ f_2 = \text{similarity}(x, l^{(2)}) \\ \vdots \end{matrix} \rightarrow \text{まとめ? } f = \begin{pmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \end{pmatrix}, \text{ ただし } f_0 = 1$$

training example  $(x^{(i)}, y^{(i)})$  に対して

$$\begin{matrix} f_1^{(i)} = \text{sim}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} = \text{sim}(x^{(i)}, l^{(2)}) \\ \vdots \\ f_i^{(i)} = \text{sim}(x^{(i)}, l^{(i)}) = 1 \quad \text{同じなので} \\ \vdots \\ f_m^{(i)} = \text{sim}(x^{(i)}, l^{(m)}) \end{matrix}$$

$\rightarrow$  まとめ?  $f^{(i)} = \begin{pmatrix} f_0^{(i)} \\ f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{pmatrix}$  ← 1 を除いて記述する場合がある。

与えられた  $x$  に対して、この  $f$  を計算する。 if  $\theta^T f \geq 0$ , predict  $y=1$

Training (学習)

$$\min_{\theta} C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T f^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$\nwarrow$   $x^{(i)}$  の代わりに  $f^{(i)}$  を使うのがポイント  
 $\nearrow$   $\theta_0$  は入っていない  
 $n=m$  だぞうた

SVM のパラメータ

$C (= \frac{1}{\lambda})$   $C \rightarrow$  大 ( $\lambda \rightarrow$  小)  $\text{bias} \rightarrow$  小  $\text{variance} \rightarrow$  大  
 ( $\because \theta_0$  以外の  $\theta_1 \sim \theta_n$  を小さくしようという影響が減るから)

$C \rightarrow$  小 ( $\lambda \rightarrow$  大)  $\text{bias} \rightarrow$  大  $\text{variance} \rightarrow$  小  
 ( $\because \theta_0$  以外の  $\theta_1 \sim \theta_n$  を小さくしようとする影響が大きくなるから)

$\alpha^2$   $\alpha^2 \rightarrow$  大 features  $f_i$  から離れていても正の値 ( $0 \cdot \theta$ ) をとる。  
 overfitting ぎみにはなるぞい。  $\nearrow$  中々りと影響が減っていく。  
 $\text{bias} \rightarrow$  大  
 $\text{variance} \rightarrow$  小

$\alpha^2 \rightarrow$  小 上の逆で  $\text{bias} \rightarrow$  小  $\text{variance} \rightarrow$  大

[講義中の Quiz]

SVM で学習させたが overfit になった。どうすればよい

- ①  $C$  を増やす ( $\lambda \rightarrow$  小)  $\text{variance} \rightarrow$  大  $\rightarrow$  overfit になる  $\times$
- ②  $C$  を減らす ( $\lambda \rightarrow$  大)  $\theta_0$  が相対的に大  $\therefore \text{bias}$  大 underfit  $\bigcirc$
- ③  $\alpha^2$  を増やす。  $f_i$  から離れる影響が中々りと underfit  $\bigcirc$
- ④  $\alpha^2$  を減らす。  $f_i$  から離れる影響が急激に  $\rightarrow$  overfit  $\times$

# SVMs in Practice

パラメータ  $C$  を選ぶ

kernel (similarity function) を選ぶ

(例) カーネルなし (linear kernel)

$$z = 1, \text{ if } \theta^T x \geq 0$$

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \geq 0$$

$n$  大  $m$  小 のとき  $x \in \mathbb{R}^{n+1}$

(例2) Gaussian kernel

$$f_i = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\alpha^2}\right), \text{ where } l^{(i)} = x^{(i)}$$

$\alpha^2$  を選ぶ必要あり  $n$  小  $m$  大 の時

Gaussian kernel を用いる前には feature scaling を行うこと

全ての similarity function が valid な kernel とは限りません。  
[ SVM パッケージの最適化が正しく動作して、発散しない ]

↓  
Mercer's Theorem  
を満たす必要はある

たくさんのおff-the-shelf kernel が利用可能である

polinomial kernel  $(x^T l + \text{定数})^{\text{次数}}$

[ 講義中の Quiz ]

$C, \alpha^2$  のような kernel のパラメータを選ぶ時、どうやって選ぶか？

① training data による結果を出すもの

② cross validation data

← ① パラメータは cross validation data  
で決定する。

③ test data

④ SVM margin が最大のもの

# SVMs in Practice

## Multi-class classification

Multi-class classification 機能を持っている SVM パッケージも多い。

一方, one-vs-all 手法を用いることもある。

## Logistic regression vs. SVMs

featureの数  $n$  と サンプルの数  $m$  を比較して

$n \rightarrow$  大, logistic regression  
 $n=10,000$  kernel なしの (linear kernel の) SVM } を使うべき  
 $m=10 \sim 1000$

$n \rightarrow$  小,  $m \rightarrow$  中, Gaussian kernel の SVM を使うべき  
 $n=1 \sim 1000, m=10 \sim 10,000$

$n \rightarrow$  小,  $m \rightarrow$  大 ほとんどの feature を追加して  
 $n=1 \sim 1000$   $m=50,000$  以上 logistic regression  
kernel なしの (linear kernel の) SVM } を使うべき

上記の全ての場合で neural network はうまく動作するか、学習に時間がかかるとかはわからない。