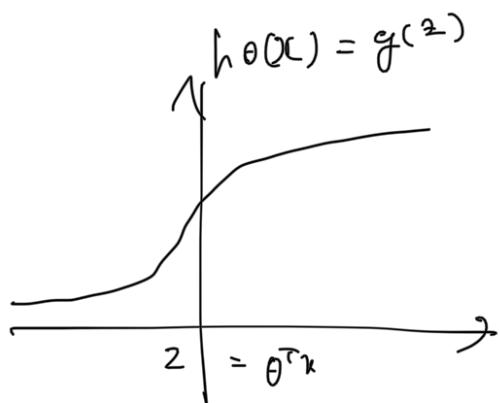


ML w7 Support Vetro Machines

Optimization Objectives.

Logistic Regression

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$



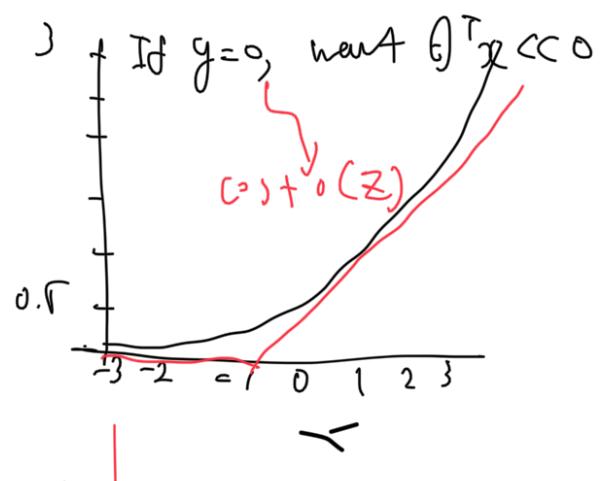
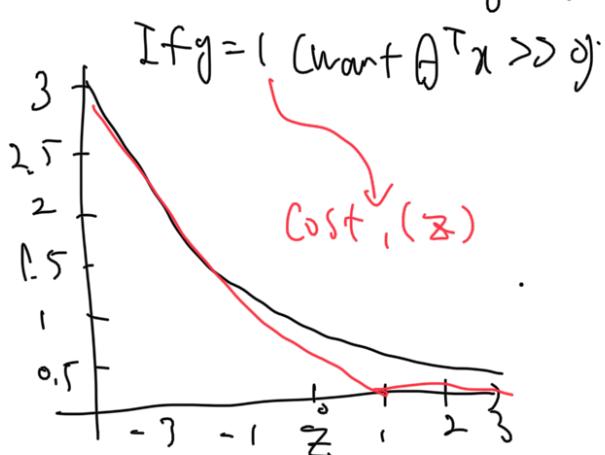
If $y=1$, we want $h_{\theta}(x) \approx 1$, $\theta^T x \gg 0$

If $y=0$, $h_{\theta}(x) \approx 0$, $\theta^T x \ll 0$

Alternative view of Logistic Regression

Cost of example: $-(y \log h_{\theta}(x) + (1-y) \log(1-h_{\theta}(x))$

$$= -y \log \frac{1}{1+e^{-\theta^T x}} - (1-y) \log \left(1 - \frac{1}{1+e^{-\theta^T x}}\right)$$



Logistic Regression

$$\min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \underbrace{\left[-1 \cdot g(h_{\theta}(x^{(i)})) \right]}_{\text{Cost} + 1 (\theta^T x^{(i)})} + (1 - y^{(i)}) \underbrace{\left[(1 - g(h_{\theta}(x^{(i)}))) \right]}_{\text{Cost} + 0 (\theta^T x^{(i)})} + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \right]$$

A B

$$\min_{\mu} (\mu - s)^2 + 1 \rightarrow \mu = s$$

$$\min_{\mu} 10(\mu - s)^2 + 10 \rightarrow \mu = s$$

$$A + \lambda B$$

$$CA + B$$

$$C = \frac{1}{\lambda}$$

$$\rightarrow \min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \max(1, (\theta^T x^{(i)}) + (1 - y^{(i)}) \max(0, (\theta^T x^{(i)})) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

Hypothesis

$$h_{\theta}(x) \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$

Large Margin Intuition

Very large number "C" and the other formula gets "0"

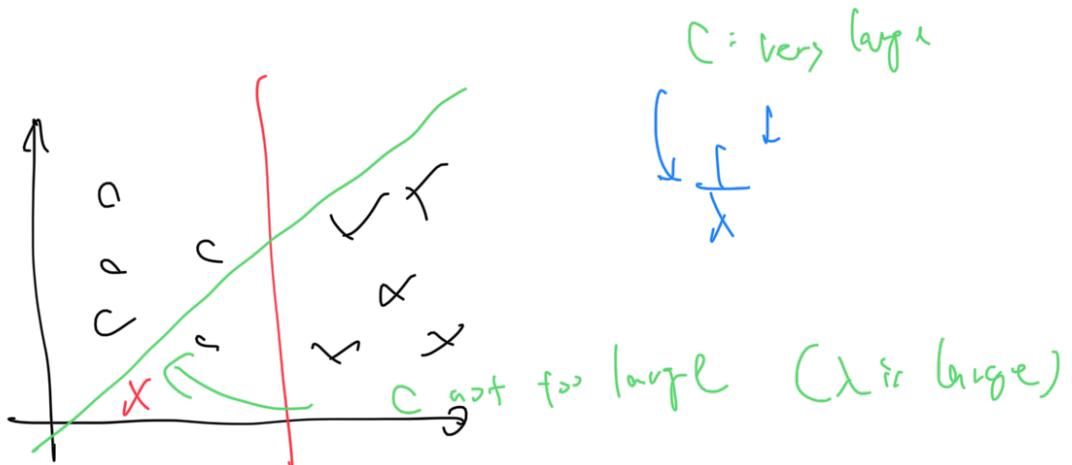
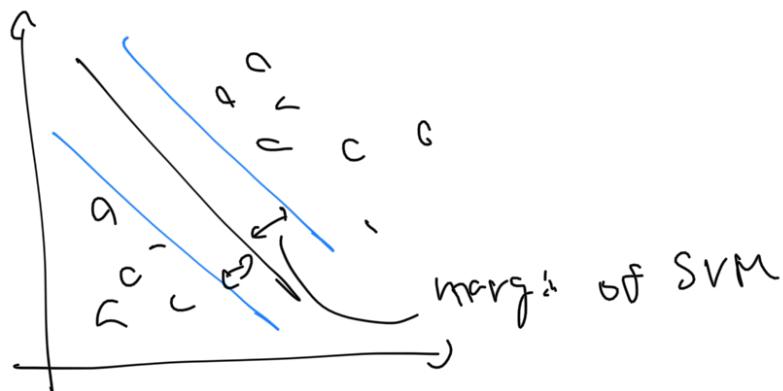
$$\min C \times 0 + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

s.t. $\theta^T x^{(i)} \geq 1$

if $y^{(i)} = 1$
if $y^{(i)} = 0$

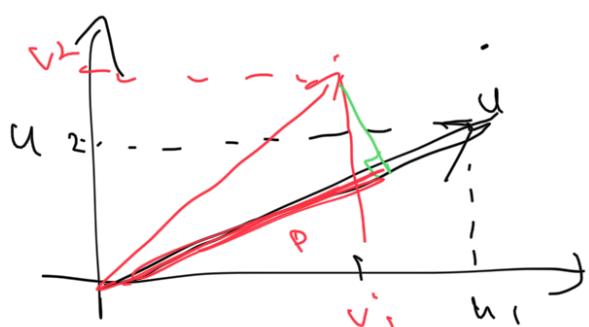
$$\theta^T x^{(i)} \leq -1 \quad \text{in } \mathcal{D}$$

SVM Decision Boundary : Linearly separable case



Mathematics Behind Large Margin Classification.

Vector Inner Product



$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$u^T v = ?$$

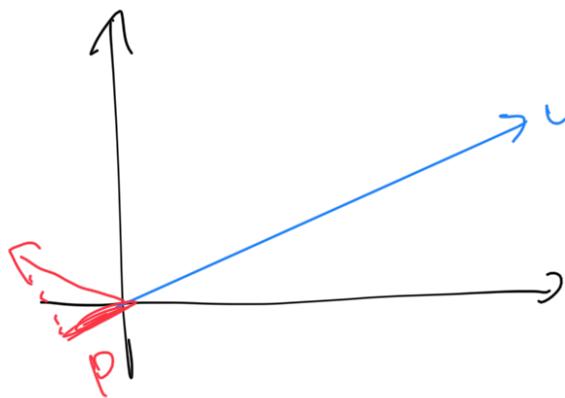
$$\|u\| = \text{length of vector } u$$

$$= \sqrt{u_1^2 + u_2^2} \in \mathbb{R}$$

$\rho = \text{length of projection of } v \text{ onto } u$

$$u^T v = p \cdot \|u\| = v^T u$$

$$= u_1 v_1 + u_2 v_2 \quad p \in \mathbb{R}$$



$$u^T v = p \cdot \|u\| \quad (p < 0)$$

SVM Decision Boundary

$$\text{min}_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} (\sqrt{\theta_1^2 + \theta_2^2})^2 = \boxed{\frac{1}{2} \|\theta\|^2}$$

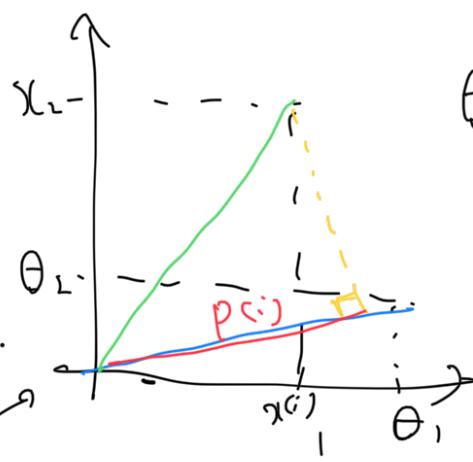
$$\text{s.t. } \theta^T x^{(i)} \geq 1 \quad \text{if } y^{(i)} = 1$$

$$\theta^T x^{(i)} \leq -1 \quad \text{if } y^{(i)} = -1$$

$$\theta^T x^{(i)} = ?$$

$$\begin{matrix} p \\ q^T \\ v \end{matrix}$$

$$\text{if } \theta_0 = 0$$

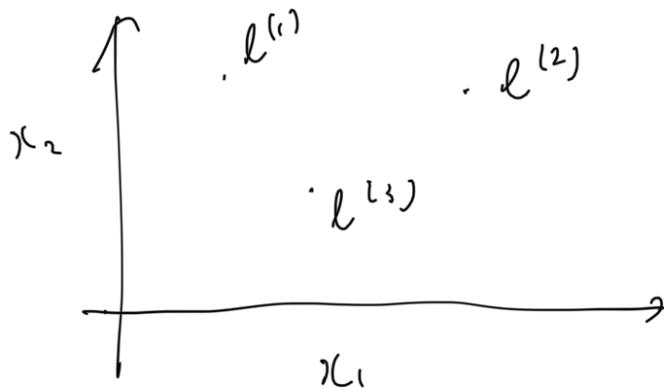


$$\begin{aligned} \theta^T x^{(i)} &= p^{(i)} \|\theta\| \\ &= \theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} \end{aligned}$$

Kernels

Non-linear Decision Boundary

Hi order polynomial dimension alternative,



Given x :

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

$$f_2 = \text{similarity}(x, l^{(2)}) = \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right)$$

kernel (Gaussian) signe.

Kernels and Similarity

$$f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\sum_{j=1}^n (x_j - l_j^{(1)})^2}{2\sigma^2}\right)$$

If $x \approx l^{(1)}$

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

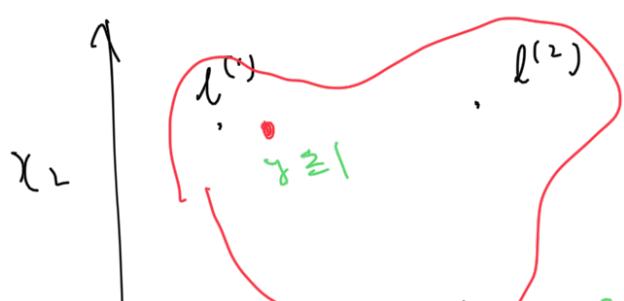
$l_1 f_1$
 $l_2 f_2$

If x far from $l^{(1)}$

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

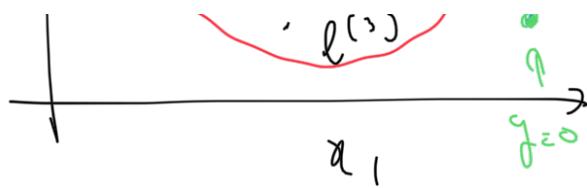
$l_3 f_3$
↑

↓ card mark



Predict "1" when

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \dots + \theta_n f_n \geq 0$$

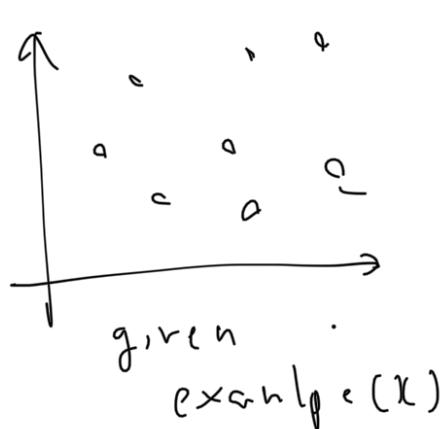


$\theta_0 = 0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0$
 $f_1 \approx x_1, f_2 \approx 0, f_3 \approx 0$.

$$\theta_0 + \theta_1 x_1 + \theta_2 x_0 + \theta_3 x_0$$

$$= -0.5 + 1 = 0.5 \geq 0$$

Kernel II.



$$(x^{(i)}, y^{(i)}), (x^{(j)}, y^{(j)})$$

...

For training example

$$\begin{aligned} x^{(i)} &\rightarrow f_1^{(i)} = \text{similarity}(x^{(i)}, l^{(i)}) \\ &f_2^{(i)} = " \\ &f_3^{(i)} = " \quad \dots \quad (x^{(i)}, l^{(i)}) \end{aligned}$$

$$\begin{aligned} x^{(i)} &\in \mathbb{R}^{n+1} \\ f^{(i)} &= \begin{bmatrix} f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} \end{aligned}$$

* SVM with kernels.

In matrix form: Given K , compute features $f \in \mathbb{R}^{m+1}$

My purpose: ...

→ Predict "y=1" if $\theta^T f \geq 0$

Training

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

$$\leftarrow \sum_j \theta_j^2 = \theta^T \theta \quad \theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_m \end{bmatrix} \quad (\text{ignore } \theta_0)$$
$$\theta^T M \theta$$

SVM parameters

$C = \left(\frac{1}{\lambda} \right)$ Large C (small λ): Low bias, High variance
small C (large λ): High bias, low variance

σ^2 large σ^2 : Feature f_i very moves smoothly
High bias, low variance

small σ^2 : Feature f_i very less smoothly
Low bias, high variance



SVM in practice,

~~for~~ liblinear, libsvm

Need specify
choice of parameter C.
choice of kernel

n : sample

m : training data

function $f = \text{kernel}(x_1, x_2)$

$$f = \exp\left(\frac{\|x_1 - x_2\|^2}{2G^2}\right)$$

→ use feature scaling before ~~using~~ using Gaussian kernel.

Other - the other kernel available:

- Polynomial kernel.

- Sigmoid kernel

- chi-square kernel.

- histogram

- intersection kernel

Multi class classification.

- use it (built-in)

- or one vs all.

Logistic regression vs SVMs.

n = number of features ($\mathbf{x} \in \mathbb{R}^{n \times 1}$)

if n is large (relative to m)

$n \geq m$, $n = [0, 1000]$, $m = [0.. \cancel{1000}]$

use logistic regression

or SVM without a kernel.

If n is small, m is intermediate,

→ use SVM with Gaussian kernel

If n is small, m is large

→ Create / add more features, then use

logistic regression or SVM without a kernel.

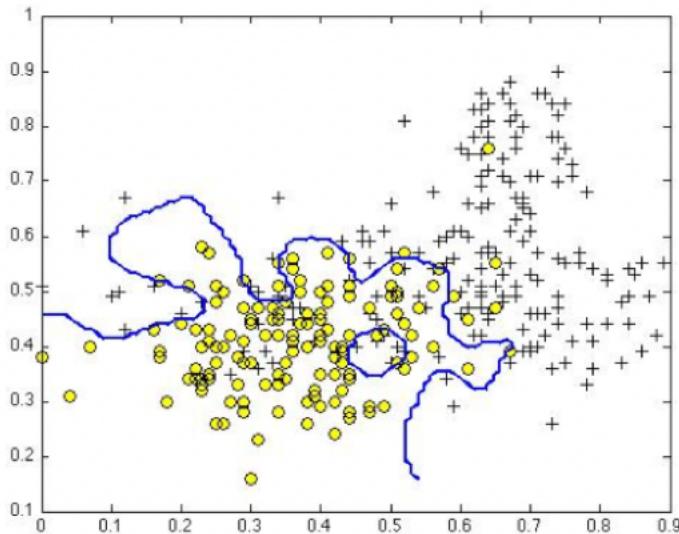
→ Neural NW works but slower.

Support Vector Machines

最新の提出物の成績評価 40%

1. Suppose you have trained an SVM classifier with a Gaussian kernel, and it learned the following decision boundary on the training set:

0 / 1点



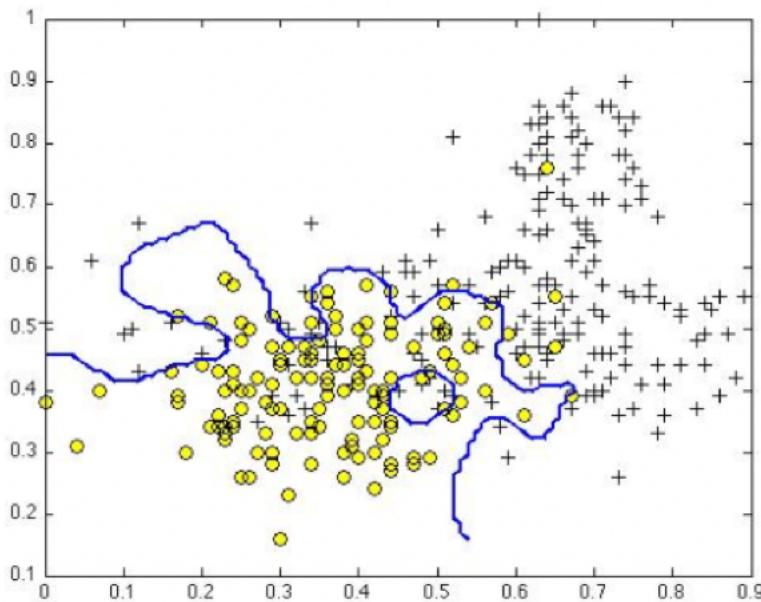
When you measure the SVM's performance on a cross validation set, it does poorly. Should you try increasing or decreasing C ? Increasing or decreasing σ^2 ?

- It would be reasonable to try **decreasing** C . It would also be reasonable to try **increasing** σ^2 .
- It would be reasonable to try **decreasing** C . It would also be reasonable to try **decreasing** σ^2 .
- It would be reasonable to try **increasing** C . It would also be reasonable to try **increasing** σ^2 .
- It would be reasonable to try **increasing** C . It would also be reasonable to try **decreasing** σ^2 .

☒ 不正解

The figure shows a decision boundary that is overfit to the training set, so we'd like to increase the bias / lower the variance of the SVM. We can do so by either decreasing (not increasing) the parameter C or increasing (not decreasing) σ^2 .

1. Suppose you have trained an SVM classifier with a Gaussian kernel, and it learned the following decision boundary on the training set:



When you measure the SVM's performance on a cross validation set, it does poorly. Should you try increasing or decreasing C ? Increasing or decreasing σ^2 ?

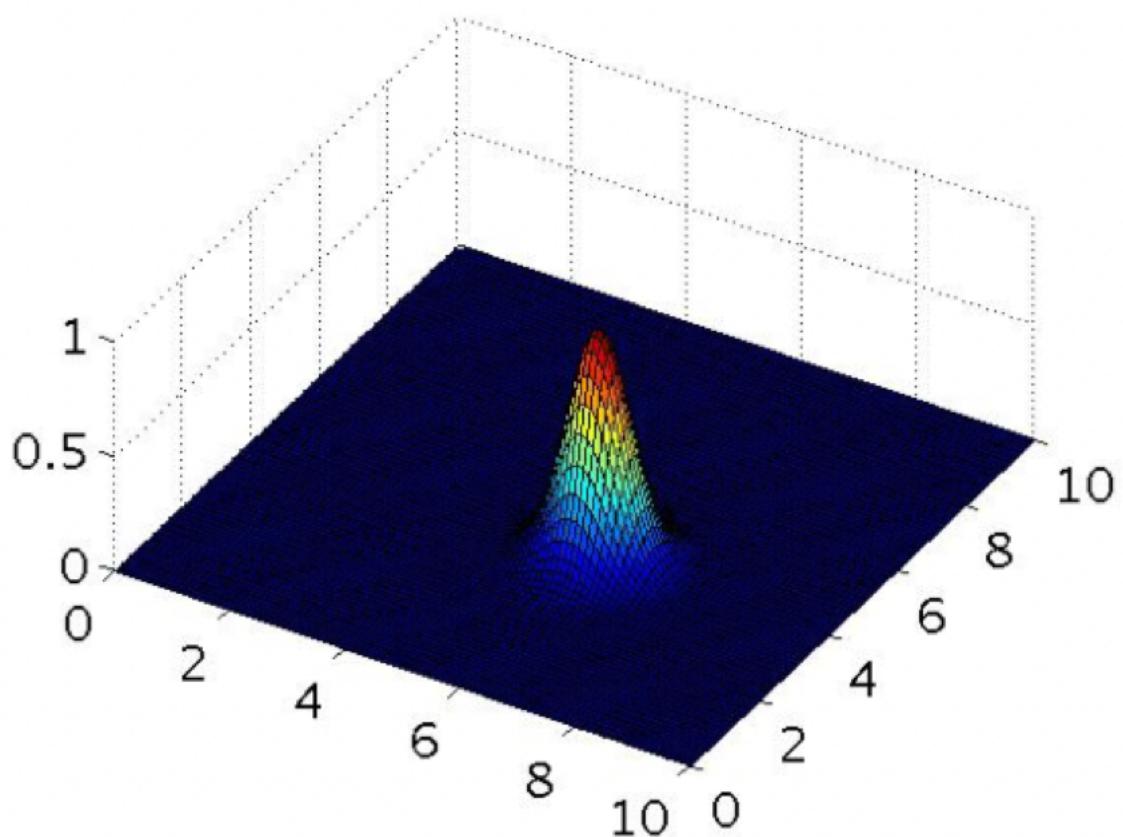
- It would be reasonable to try **decreasing** C . It would also be reasonable to try **increasing** σ^2 .
- It would be reasonable to try **decreasing** C . It would also be reasonable to try **decreasing** σ^2 .
- It would be reasonable to try **increasing** C . It would also be reasonable to try **decreasing** σ^2 .
- It would be reasonable to try **increasing** C . It would also be reasonable to try **increasing** σ^2 .



The figure shows a decision boundary that is overfit to the training set, so we'd like to increase the bias / lower the variance of the SVM. We can do so by either decreasing the parameter C or increasing σ^2 .

①

Figure 2.

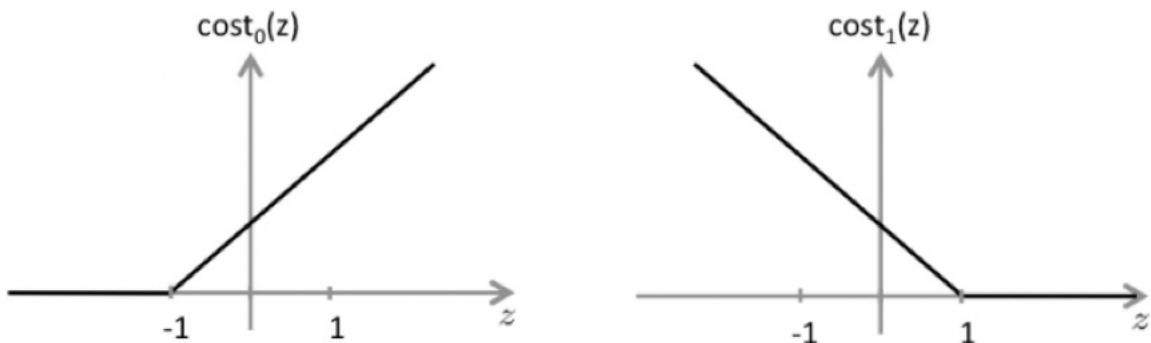


3. The SVM solves

0 / 1点

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

where the functions $\text{cost}_0(z)$ and $\text{cost}_1(z)$ look like this:



The first term in the objective is:

$$C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}).$$

This first term will be zero if two of the following four conditions hold true. Which are the two conditions that would guarantee that this term equals zero?

- For every example with $y^{(i)} = 0$, we have that $\theta^T x^{(i)} \leq -1$.

Ⓐ 正解

For examples with $y^{(i)} = 0$, only the $\text{cost}_0(\theta^T x^{(i)})$ term is present. As you can see in the graph, this will be zero for all inputs less than or equal to -1.

- For every example with $y^{(i)} = 1$, we have that $\theta^T x^{(i)} \geq 0$.
- For every example with $y^{(i)} = 1$, we have that $\theta^T x^{(i)} \geq 1$.
- For every example with $y^{(i)} = 0$, we have that $\theta^T x^{(i)} \leq 0$.

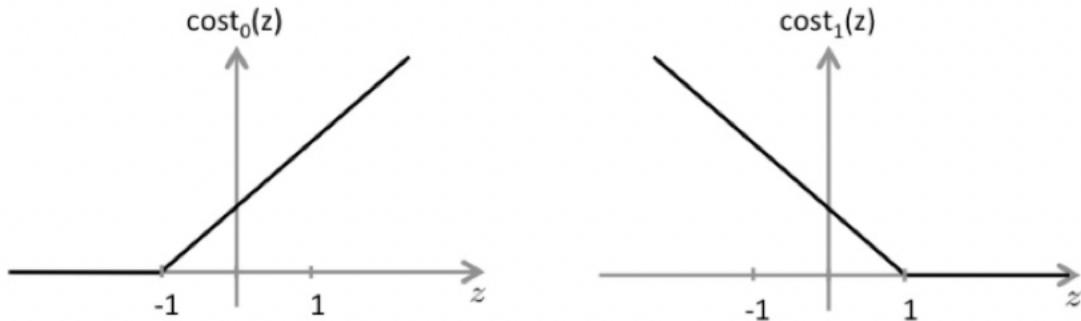
正しい回答をすべて選択しました

3. The SVM solves

0 / 1点

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

where the functions $\text{cost}_0(z)$ and $\text{cost}_1(z)$ look like this:



The first term in the objective is:

$$C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}).$$

This first term will be zero if two of the following four conditions hold true. Which are the two conditions that would guarantee that this term equals zero?

- For every example with $y^{(i)} = 0$, we have that $\theta^T x^{(i)} \leq -1$.



正解

For examples with $y^{(i)} = 0$, only the $\text{cost}_0(\theta^T x^{(i)})$ term is present. As you can see in the graph, this will be zero for all inputs less than or equal to -1.

- For every example with $y^{(i)} = 0$, we have that $\theta^T x^{(i)} \leq 0$.



× これを選択しないでください

$\text{cost}_0(\theta^T x^{(i)})$ is still non-zero for inputs between -1 and 0, so being less than or equal to 0 is insufficient.

- For every example with $y^{(i)} = 1$, we have that $\theta^T x^{(i)} \geq 0$.

- For every example with $y^{(i)} = 1$, we have that $\theta^T x^{(i)} \geq 1$.

4. Suppose you have a dataset with $n = 10$ features and $m = 5000$ examples.

0 / 1点

After training your logistic regression classifier with gradient descent, you find that it has underfit the training set and does not achieve the desired performance on the training or cross validation sets.

Which of the following might be promising steps to take? Check all that apply.

- Use an SVM with a Gaussian Kernel.

✓ 正解

By using a Gaussian kernel, your model will have greater complexity and can avoid underfitting the data.

- Increase the regularization parameter λ .

✗ これを選択しないでください

You are already underfitting the data, and increasing the regularization parameter only makes underfitting stronger.

- Use an SVM with a linear kernel, without introducing new features.

- Create / add new polynomial features.

4. Suppose you have a dataset with $n = 10$ features and $m = 5000$ examples.

After training your logistic regression classifier with gradient descent, you find that it has underfit the training set and does not achieve the desired performance on the training or cross validation sets.

Which of the following might be promising steps to take? Check all that apply.

Increase the regularization parameter λ .

Use an SVM with a Gaussian Kernel.



By using a Gaussian kernel, your model will have greater complexity and can avoid underfitting the data.

Create / add new polynomial features.



When you add more features, you increase the variance of your model, reducing the chances of underfitting.

Use an SVM with a linear kernel, without introducing new features.

5. Which of the following statements are true? Check all that apply.

The maximum value of the Gaussian kernel (i.e., $\text{sim}(x, l^{(1)})$) is 1.

 正解

When $x = l^{(1)}$, the Gaussian kernel has value $\exp(0) = 1$, and it is less than 1 otherwise.

If the data are linearly separable, an SVM using a linear kernel will

return the same parameters θ regardless of the chosen value of

C (i.e., the resulting value of θ does not depend on C).

Suppose you are using SVMs to do multi-class classification and

would like to use the one-vs-all approach. If you have K different

classes, you will train $K - 1$ different SVMs.

It is important to perform feature normalization before using the Gaussian kernel.

 正解

The similarity measure used by the Gaussian kernel expects that the data lie in approximately the same range.