

Large Margin Classification

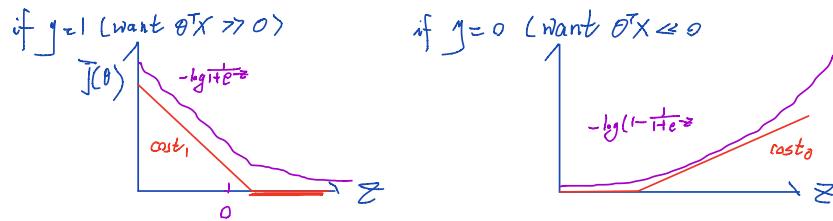
Optimization Objective

Alternative view of logistic regression

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

$$J(\theta) = -y \log h_{\theta}(x) + (1-y) \log(1-h_{\theta}(x))$$

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



Support Vector machine (SVM)

logistic regression $\min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} (-\log h_{\theta}(x^{(i)})) + (1-y^{(i)}) (\underline{-\log(1-h_{\theta}(x^{(i)}))}) \right] + \frac{\lambda}{2} \sum_{j=0}^n \theta_j^2$

Support vector machine $\min_{\theta} \frac{1}{m} \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] + \frac{\lambda}{2} \sum_{j=0}^n \theta_j^2$

In logistic regression $A + \lambda B$

In Support Vector machine $CA + B \quad (C = \frac{1}{\lambda})$

Cost function for SVM:

$$\min_{\theta} \left[\sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] \right] + \frac{\lambda}{2} \sum_{j=0}^n \theta_j^2$$

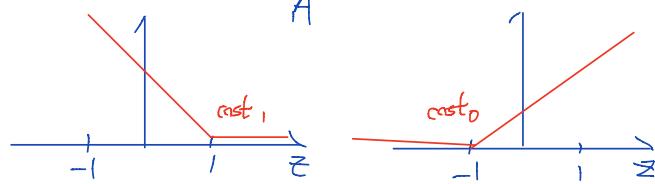
Hypothesis of SVM

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

... | if else

Large Margin intuition

$$\text{SVM: } \min_C \sum_{i=1}^m [\underbrace{y^{(i)} \text{cost}_1(\theta^T x^{(i)})}_{A''} + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=1}^n Q_j^2$$



if $y=1$, we want $\theta^T x \geq 1$ (not just > 0)

if $y=0$, we want $\theta^T x \leq -1$ (not just < 0)

Suppose C very large, $\min_C C A + B \Rightarrow A \rightarrow 0$

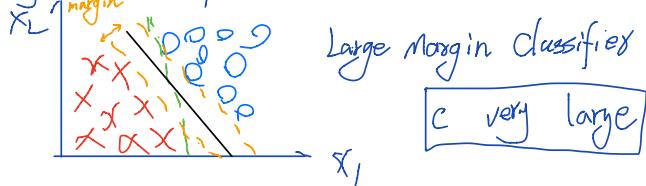
whenever $y^{(i)}=1$:

$$\theta^T x \geq 1, \min C \times 0 + \frac{1}{2} \sum_{j=1}^n Q_j^2$$

whenever $y^{(i)}=0$

$$\theta^T x \leq -1, \min C \times 0 + \frac{1}{2} \sum_{j=1}^n Q_j^2$$

SVM Decision Boundary: Linearly separable case



SVM is a large margin classifier

Mathematics behind Large Margin Classification

SVM decision boundary when C is very large:

(for simplicity, assume $n=2$, $\theta_0=0$)

$$\min \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} \|\theta\|^2 \quad \theta = [\theta_1 \quad \theta_2]$$

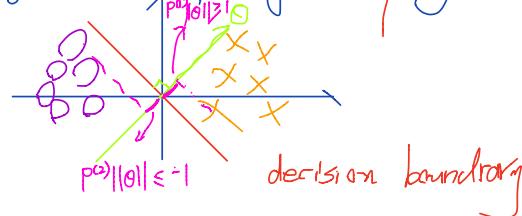
if $y^{(i)} = 1, \theta^T x^{(i)} \geq 1 \Rightarrow p^{(i)} \cdot \|\theta\| \geq 1$

if $y^{(i)} = 0, \theta^T x^{(i)} \leq -1 \Rightarrow p^{(i)} \cdot \|\theta\| \leq -1$

where $p^{(i)}$ is $x^{(i)}$ projection on to θ

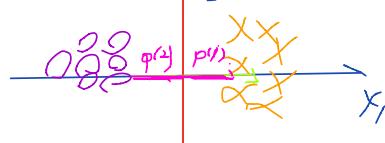
Now prove that SVM with large C will end up with large decision margin (large margin classification)

1) SVM (C large) will not do following small margin:



projection $p^{(i)}$ is small. means $\|\theta\|$ has to be large to make $p^{(i)} \|\theta\| \geq 1$ and $p^{(i)} \|\theta\| \leq -1$, but cost function want $\frac{1}{2} \|\theta\|^2$ to be small, so SVM (large C) will not end up with this situation

2) SVM (large C) will try to do following



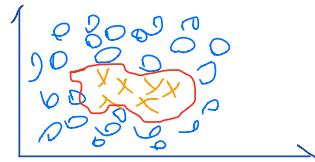
In this case $p^{(i)}$ is larger, to fulfil $p^{(i)} \|\theta\| \geq 1$ and $p^{(i)} \|\theta\| \leq -1$, $\|\theta\|$ doesn't need to be too large, which agrees the minimization of cost function $\min \frac{1}{2} \|\theta\|^2$

thus SVM (with large C) is a Large Margin Classification

Assumption $\theta_0=0$ make θ go through origin of wbar plot

Kernels.

Nonlinear Decision Boundary



predict $y=1$ if $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \dots > 0$

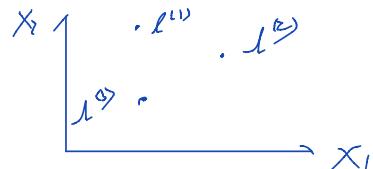
$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta_0 + \theta_1 x_1 + \dots > 0 \\ 0 & \text{else} \end{cases}$$

can rewrite: $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \dots$

where: $f_1 = x_1, f_2 = x_2, f_3 = x_1 x_2, \dots$

Q: How to optimize the choice of f_1, f_2, f_3, \dots ?

Kernel



Pick landmarks $l^{(1)}, l^{(2)}, l^{(3)}$, given x , compute new feature depend on proximity to landmarks $l^{(1)}, l^{(2)}, l^{(3)}$

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

$$f_2 = \text{similarity}(x, l^{(2)}) = \exp(\dots)$$

$$f_3 = \dots$$

kernel $k(x, l^{(1)})$ Gaussian kernel

kernels and similarity

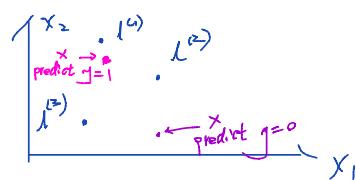
$$f_i = \text{similarity}(x, l^{(i)}) = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\sum_{j=1}^n (x_j - l_j^{(i)})^2}{2\sigma^2}\right)$$

if $x \approx l^{(i)}$: $f_i \approx 1$

if x far from $l^{(i)}$, $f_i \approx 0$

each of landmarks defines a new feature

Example:



predict 1 when $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \theta_4 f_4 \geq 0$

say $\theta_0 = -0.5$, $\theta_1 = 1$, $\theta_2 = 1$, $\theta_3 = 0$

case 1: $f_1 \approx 1$, $f_2 \approx 0$, $f_3 \approx 0$

$$\theta_0 + \theta_1 x_1 + \theta_2 x_0 + \theta_3 x_0 = -0.5 + 1 = 0.5 > 0 \quad \text{predict 1}$$

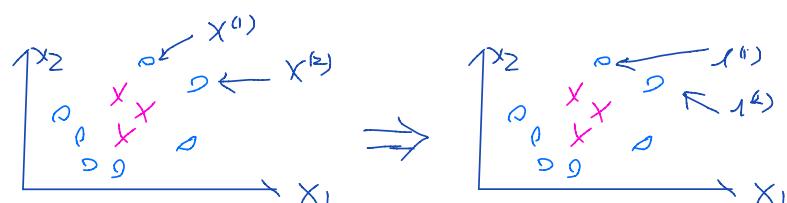
case 2: $f_1, f_2, f_3 \approx 0$

$$\theta_0 + \theta_1 x_0 + \theta_2 x_0 + \theta_3 x_0 = -0.5 < 0 \quad \text{predict 0}$$

Understand: for points close to landmarks, predict 1

for points far away from landmarks, predict 0

choosing the landmarks



choose landmark exactly where the training data is
in training examples \Rightarrow m landmarks

Given $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$
choose $\ell^{(1)} = x^{(1)}, \dots, \ell^{(m)} = x^{(m)}$

Given example x :

$$f_1 = \text{Similarity}(x, \ell^{(1)})$$

$$f_2 = \dots$$

$$f_m = \dots$$

$$f = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix} \rightarrow = 1$$

For training example $(x^{(i)}, y^{(i)})$

$$f_0^{(i)} = \text{Sim}(x^{(i)}, \ell^{(0)})$$

$$f_1^{(i)} = \text{Sim}(x^{(i)}, \ell^{(1)})$$

$$\vdots \quad \vdots \quad \leftarrow f_i^{(i)} = \text{Sim}(x^{(i)}, \ell^{(i)}) = 1$$

$$f_m^{(i)} = \text{Sim}(x^{(i)}, \ell^{(m)})$$

$$\text{then: } x \in \mathbb{R}^{n+1} \rightarrow f^{(i)} = \begin{bmatrix} f_0^{(i)} = 1 \\ f_1^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} \text{ new feature function}$$

SVM with kernels ($n=m$ or $n=m+1$)

Hypothesis: Given x . compute features $f \in \mathbb{R}^{n+1}$

predict " $y=1$ " if $\theta^T f \geq 0$ $\theta^T f = \theta_0 f_0 + \theta_1 f_1 + \dots + \theta_m f_m$ $\theta \in \mathbb{R}^{n+1}$

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

In real software/packages, it's actually $\frac{1}{2} \sum_{j=1}^{n+1} \theta_j^2$ instead of $\|\theta\|^2$
where M is some kernel related matrix, in order to let the algorithm scale to much large datasets and computes faster

* SVM with kernel usually runs slow

SVM parameters:

$C (= \frac{1}{\lambda})$: Large C : lower bias, high variance
Small C : higher bias, low variance

σ^2 : large σ^2 : features f_i vary more smoothly
higher bias, lower variance
Small σ^2 : feature f_i vary less smoothly
lower bias, higher variance

Using An SVM

svm software package: liblinear, libsvm

Need to specify:

- choice of parameter C
- choice of kernel

e.g. No kernel ("linear kernel")

predict " $y=1$ " if $\theta^T X \geq 0$

- good for n (# of features) large, m (# of examples) small
to avoid overfitting

e.g. Gaussian kernel

need to choose σ^2

* good for n small, m large

* Note: Do perform feature scaling before using Gaussian kernel, reason is

$$\|x - \mu\|^2 = (x_1 - \mu_1)^2 + (x_2 - \mu_2)^2 + \dots + (x_n - \mu_n)^2$$

$$\sim 10^3 \quad \sim 1$$

then x_2 doesn't show much influence

other choice of kernel:

Note: Not all similarity functions make valid kernels

Need to satisfy "Mercer's Theorem" to make sure SVM packages' optimizations run correctly and do not diverge.

Many off-the-shelf kernels available:

- polynomial kernel: $(x^T x + \text{const})^{\text{degree}} (=2, 3, \dots)$
- string kernel
- chi-square kernel
- histogram intersection kernel

Multi-class classification

Many SVM packages already have built-in multi-class classification functionality

Otherwise, use one-vs-all method. Train K SVMs, one to distinguish $y = i$ from the rest, for $i = 1, 2, \dots, K$, get $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)}$, pick class i with largest $(\theta^{(i)})^T x$

Logistic regression vs. SVMs

How to make choice in between these two algorithms?

n : number of features ($x \in \mathbb{R}^{n+1}$)

m: number of training examples

- o n is large relative to m (e.g. $n \gg m$, $n=10,000$, $n=10, \dots, 100$)
 - ↳ Use regression, or SVM w/o a kernel ("linear kernel")
- o n is small, m is intermediate (e.g. $n=1-1000$, $m=10-10,000$)
 - ↳ Use SVM w/ Gaussian kernel
- o n is small, m is large ($n=1 \sim 1000$, $M=50,000 + >$)
 - ↳ Create/add more features, then use logistic regression or SVM w/o a kernel ("linear kernel")
- o neural network likely to work well for most of these settings, but may be slower to train

* SVM is a convex problem, can always find ~~the~~ global min