

# Using an SVM

*SVM in Practice*

Support Vector Machines

Use SVM software package (e.g. liblinear, libsvm, ...) to solve for parameters  $\theta$ .

Use SVM software package (e.g. liblinear, libsvm, ...) to solve for parameters  $\theta$ .

Need to specify:

- Choice of parameter  $C$ .
- Choice of kernel (similarity function):

Use SVM software package (e.g. liblinear, libsvm, ...) to solve for parameters  $\theta$ .

Need to specify:

- Choice of parameter  $C$ .
- Choice of kernel (similarity function):

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

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

Use SVM software package (e.g. liblinear, libsvm, ...) to solve for parameters  $\theta$ .



Need to specify:

→ Choice of parameter  $C$ .

Choice of kernel (similarity function):

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

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

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

$n$  large,  $m$  small

Use SVM software package (e.g. liblinear, libsvm, ...) to solve for parameters  $\theta$ .



Need to specify:

→ Choice of parameter  $C$ .

Choice of kernel (similarity function):

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

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

$$\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n \geq 0 \quad \rightarrow \quad \underline{n} \text{ large}, \quad \underline{m} \text{ small} \quad \underline{x \in \mathbb{R}^{n+1}}$$

Gaussian kernel:

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

Need to choose  $\sigma^2$ .

Use SVM software package (e.g. liblinear, libsvm, ...) to solve for parameters  $\theta$ .



Need to specify:

→ Choice of parameter C.

Choice of kernel (similarity function):

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

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

$$\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n \geq 0 \quad \rightarrow \underline{n} \text{ large, } \underline{m} \text{ small} \quad \underline{x \in \mathbb{R}^{n+1}}$$

Gaussian kernel:

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

Need to choose  $\sigma^2$



$x \in \mathbb{R}^n$ ,  $n$  small  
and/or  $m$  large



Kernel (similarity) functions:

```
function f = kernel(x1,x2)
```

$$f = \exp \left( -\frac{\| \mathbf{x1} - \mathbf{x2} \|^2}{2\sigma^2} \right)$$

```
return
```



Kernel (similarity) functions:

function  $f = \text{kernel}(\underline{x1}, \underline{x2})$

$$f = \exp\left(-\frac{\|\underline{x1} - \underline{x2}\|^2}{2\sigma^2}\right)$$

return

$x \rightarrow$

$$\begin{matrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{matrix}$$

Kernel (similarity) functions:

function  $f = \text{kernel}(\underline{x1}, \underline{x2})$

$$f = \exp\left(-\frac{\|\underline{x1} - \underline{x2}\|^2}{2\sigma^2}\right)$$

return

$\leftarrow$

$x \rightarrow$

$f_1$   
 $f_2$   
 $\vdots$   
 $f_m$

→ Note: Do perform feature scaling before using the Gaussian kernel.

## Kernel (similarity) functions:

function  $f = \text{kernel}(\mathbf{x}_1, \mathbf{x}_2)$

$$f = \exp\left(-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|^2}{2\sigma^2}\right)$$

return

$x \rightarrow \begin{matrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{matrix}$

→ Note: Do perform feature scaling before using the Gaussian kernel.

$$\rightarrow \|\mathbf{x} - \mathbf{l}\|^2$$

$$\mathbf{v} = \mathbf{x} - \mathbf{l}$$

$$\|\mathbf{v}\|^2 = v_1^2 + v_2^2 + \dots + v_n^2$$

$$= (x_1 - l_1)^2 + (x_2 - l_2)^2 + \dots + (x_n - l_n)^2$$

1000 feet<sup>2</sup>

1-5 bedrooms

## Other choices of kernel

Note: Not all similarity functions  $\text{similarity}(x, l)$  make valid kernels.  
(Need to satisfy technical condition called “Mercer’s Theorem” to make sure SVM packages’ optimizations run correctly, and do not diverge).

## Other choices of kernel

Note: Not all similarity functions  $\text{similarity}(x, l)$  make valid kernels.

→ (Need to satisfy technical condition called “Mercer’s Theorem” to make sure SVM packages’ optimizations run correctly, and do not diverge).

Many off-the-shelf kernels available:

- Polynomial kernel:

$$k(x, l) = (x^T l)^3, \quad (x^T l)^2, \quad (x^T l + 1)^3, \quad (x^T l + 5)^4, \quad (x^T l + \text{constant})^{\text{degree}}$$

Handwritten annotations: Red arrows point to the exponents 3, 2, 3, 4, and the word "degree". A red arrow also points to the constant term in the last expression.



## Other choices of kernel

Note: Not all similarity functions  $\text{similarity}(x, l)$  make valid kernels.

→ (Need to satisfy technical condition called “Mercer’s Theorem” to make sure SVM packages’ optimizations run correctly, and do not diverge).

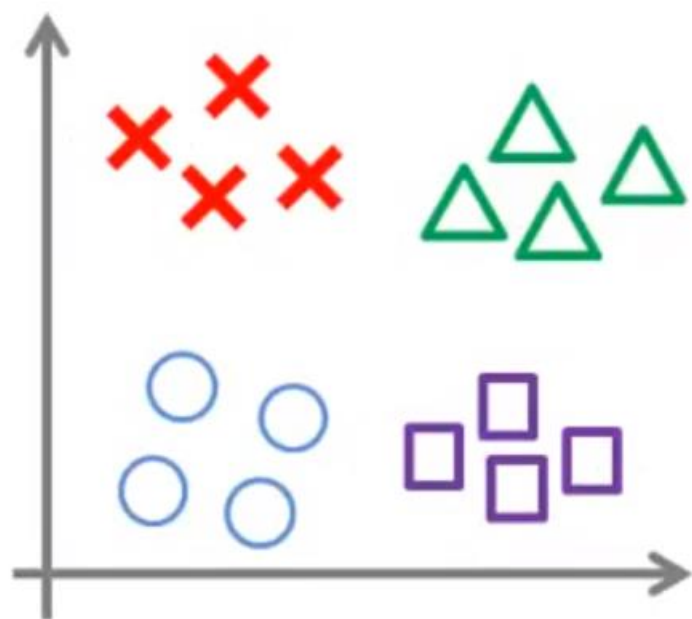
Many off-the-shelf kernels available:

- Polynomial kernel:

$$k(x, l) = (x^T l)^3, \quad (x^T l)^2, \quad (x^T l + 1)^3, \quad (x^T l + 5)^4, \quad (x^T l + \text{constant})^{\text{degree}}$$

- More esoteric: String kernel, chi-square kernel, histogram intersection kernel, ...

## Multi-class classification



$$y \in \{1, 2, 3, \dots, K\}$$

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

$n$  = number of features ( $x \in \mathbb{R}^{n+1}$ ),  $m$  = number of training examples

If  $n$  is large (relative to  $m$ ):

Use logistic regression, or SVM without a kernel (“linear kernel”)



## Logistic regression vs. SVMs

$n$  = number of features ( $x \in \mathbb{R}^{n+1}$ ),  $m$  = number of training examples

→ If  $n$  is large (relative to  $m$ ): (e.g.  $n \geq m$ ,  $n = 10,000$ ,  $m = 10 \dots 1000$ )

Use logistic regression, or SVM without a kernel ("linear kernel")

## Logistic regression vs. SVMs

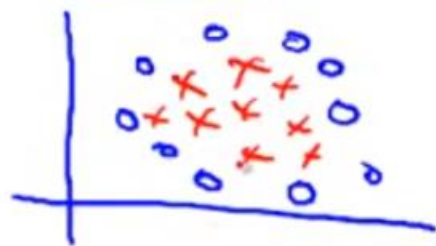
$n$  = number of features ( $x \in \mathbb{R}^{n+1}$ ),  $m$  = number of training examples

→ If  $n$  is large (relative to  $m$ ): (e.g.  $n \geq m$ ,  $n = \underline{10,000}$ ,  $m = \underline{10} \dots \underline{1000}$ )

→ Use logistic regression, or SVM without a kernel ("linear kernel")

→ If  $n$  is small,  $m$  is intermediate: ( $n = 1 - 1000$ ,  $m = 10 - \underline{10,000}$ )

→ Use SVM with Gaussian kernel



## Logistic regression vs. SVMs

$n$  = number of features ( $x \in \mathbb{R}^{n+1}$ ),  $m$  = number of training examples

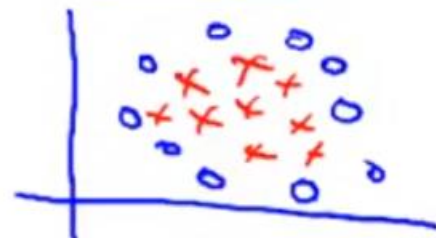
→ If  $n$  is large (relative to  $m$ ): (e.g.  $n \geq m$ ,  $n = \underline{10,000}$ ,  $m = \underline{10} \dots \underline{1000}$ )

→ Use logistic regression, or SVM without a kernel ("linear kernel")

→ If  $n$  is small,  $m$  is intermediate: ( $n = 1 - 1000$ ,  $m = 10 - \underline{10,000}$ )

→ Use SVM with Gaussian kernel

If  $n$  is small,  $m$  is large: ( $n = 1 - 1000$ ,  $m = 50,000 +$ )



Create/add more features, then use logistic regression or SVM without a kernel

Packages are good but they struggle if there are too many data.

## Logistic regression vs. SVMs

$n$  = number of features ( $x \in \mathbb{R}^{n+1}$ ),  $m$  = number of training examples

→ If  $n$  is large (relative to  $m$ ): (e.g.  $n \geq m$ ,  $n = \underline{10,000}$ ,  $m = \underline{10} \dots \underline{1000}$ )

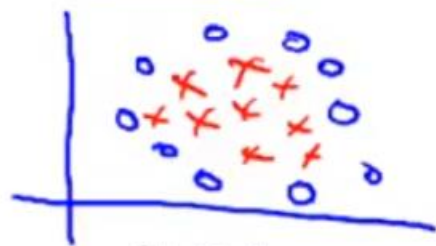
→ Use logistic regression, or SVM without a kernel ("linear kernel")

→ If  $n$  is small,  $m$  is intermediate: ( $n = \underline{1-1000}$ ,  $m = \underline{10-10,000}$ ) ←

→ Use SVM with Gaussian kernel

If  $n$  is small,  $m$  is large: ( $n = \underline{1-1000}$ ,  $m = \underline{50,000+}$ )

→ Create/add more features, then use logistic regression or SVM without a kernel



→ Neural network likely to work well for most of these settings, but may be slower to train.