#### **Support Vector Machines** INFO-4604, Applied Machine Learning University of Colorado Boulder

**September 27, 2018** 

Prof. Michael Paul

# Today

Two important concepts:

- Margins
- Kernels



#### **Linear Predictions**

Perceptron:

$$f(\mathbf{x}) = \begin{bmatrix} 1, & \mathbf{w}^{\mathsf{T}}\mathbf{x} \ge 0\\ -1, & \mathbf{w}^{\mathsf{T}}\mathbf{x} < 0 \end{bmatrix}$$

SVM:

$$f(\mathbf{x}) = \begin{bmatrix} 1, & \mathbf{w}^{\mathsf{T}}\mathbf{x} \ge 1 \\ -1, & \mathbf{w}^{\mathsf{T}}\mathbf{x} \le -1 \end{bmatrix}$$

Two different boundaries for positive vs negative



The **margin** is the distance between the two boundaries.

The **support vectors** are the instances at the boundaries (when  $\mathbf{w}^T \mathbf{x} = 1$  or -1)

• Or within the boundaries, if not linearly separable

The goal of SVMs is to learn the boundaries to make the margin as large as possible (while still correctly classifying the instances)

• maximum margin classification

The size of the margin is: 2 / IIwII

- Recall: IIwII is the L2 norm of the weight vector
- Smaller weights  $\rightarrow$  larger margin

#### Learning goal:

- Maximize 2 / IIwII, subject to the constraints that all instances are correctly classified
- Turn it into minimization problem by taking the inverse: ½ IIwII
- Can also square the L2 norm (makes the calculus easier), just like with L2 regularization: <sup>1</sup>/<sub>2</sub> IIwII<sup>2</sup>

The size of the margin is: 2 / IIwII

- Recall: IIwII is the L2 norm of the weight vector
- Smaller weights  $\rightarrow$  larger margin

#### Learning goal:

- Minimize:  $\frac{1}{2} \| \boldsymbol{w} \|^2$
- Subject to the constraints:

Only possible to satisfy these constraints if the instances are linearly separable!

$$y^{(i)}\left(w_0 + \boldsymbol{w}^T \boldsymbol{x}^{(i)}\right) \ge 1 \; \forall_i$$

In the general case, SVM uses this loss function:

$$L_{i}(\mathbf{w}; \mathbf{x}_{i}) = \begin{bmatrix} 0, & y_{i} (\mathbf{w}^{\mathsf{T}} \mathbf{x}_{i}) \geq \\ -y_{i} (\mathbf{w}^{\mathsf{T}} \mathbf{x}_{i}), & \text{otherwise} \end{bmatrix}$$

Same as perceptron, but  $y_i (\mathbf{w}^T \mathbf{x}_i) \ge 1$  instead of  $y_i (\mathbf{w}^T \mathbf{x}_i) \ge 0$ 



In the general case, SVM uses this loss function:

$$L_{i}(\mathbf{w}; \mathbf{x}_{i}) = \begin{bmatrix} 0, & y_{i} (\mathbf{w}^{\mathsf{T}} \mathbf{x}_{i}) \ge 1 \\ -y_{i} (\mathbf{w}^{\mathsf{T}} \mathbf{x}_{i}), & \text{otherwise} \end{bmatrix}$$

The learning goal of SVMs when the data are not linearly separable is to minimize:



SVMs also use L2 regularization

• *C* is like  $\lambda$  from before, but larger *C*  $\rightarrow$  lower loss

Like perceptron, the SVM function can be minimized using stochastic (sub)gradient descent

• With sklearn's SGDClassifier class, SVM can be implemented by setting loss='hinge'

Other implementations (usually using different optimization algorithms than SGD)

- Liblinear and LIBSVM (both used by sklearn)
- SVM-light

# Large Margins: Summary



# Large Margins: Summary

Classifiers with large margins are more likely to have better generalization, less overfitting

• Hyperparameter *C* controls the tradeoff between margin size and classification error

The large margin principle is another justification for L2 regularization that you saw earlier

• Since the size of the margin is inversely proportional to the L2 norm of the weight vector

## Kernel Trick

It turns out that the optimal solution for **w** is equivalent to:



Combination of each training instance's feature vector, weighted by  $\boldsymbol{\alpha}$ 

So in the loss function and prediction functions, we can replace  $\mathbf{w}^{\mathsf{T}}\mathbf{x}$  with  $\sum_{i} \alpha_{i} \mathbf{x}_{i}^{\mathsf{T}}\mathbf{x}$ 

 $\alpha_i$  is only nonzero for support vectors

• This summation can therefore skip over all other instances, making this calculation more efficient.

### Kernel Trick

In the loss function and prediction functions, we can replace  $\mathbf{w}^{\mathsf{T}}\mathbf{x}$  with  $\sum_{i} \alpha_{i} \mathbf{x}_{i}^{\mathsf{T}}\mathbf{x}$ 

Now this looks similar to weighted nearest neighbor classification, where the "similarity" between an instance **x** and another instance **x**<sub>i</sub> is  $\mathbf{x_i}^{\mathsf{T}}\mathbf{x}$  and this is additionally weighted by  $\alpha_i$ 

Learning goal is now to learn  $\alpha$  instead of  $\boldsymbol{w}$ 

• How? More complex than before...

### **Kernel Functions**

Loosely, a kernel function is a similarity function between two instances

General kernel trick: Replace  $\mathbf{w}^{\mathsf{T}}\mathbf{x}$  with  $\sum_{i} \alpha_{i} k(\mathbf{x}_{i}, \mathbf{x})$ 

The **linear kernel** function for an SVM is:  $k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$ 

### **Kernel Functions**

What happens if we define the kernel function in some other way?

Then it won't be true that  $\sum_{i} \alpha_{i} k(\mathbf{x}_{i}, \mathbf{x}) = \mathbf{w}^{\mathsf{T}} \mathbf{x}$ 

But: kernels can be defined so that,

$$\sum_{\mathbf{i}} \alpha_{\mathbf{i}} \, \mathsf{k}(\mathbf{x}_{\mathbf{i}}, \, \mathbf{x}) = \mathbf{w}^{\mathsf{T}} \boldsymbol{\phi}(\mathbf{x}),$$

where  $\phi(\mathbf{x})$  is some other feature representation.

#### Kernel Functions: Polynomial

#### A **polynomial kernel** function is defined as: $k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j + c)^d$

If d=2 (quadratic kernel), then it turns out that  $\sum_{i} \alpha_{i} k(\mathbf{x}_{i}, \mathbf{x}) = \mathbf{w}^{T} \boldsymbol{\phi}(\mathbf{x})$ 

where

$$arphi(x)=\langle x_n^2,\ldots,x_1^2,\sqrt{2}x_nx_{n-1},\ldots,\sqrt{2}x_nx_1,\sqrt{2}x_{n-1}x_{n-2},$$

$$\ldots,\sqrt{2}x_{n-1}x_1,\ldots,\sqrt{2}x_2x_1,\sqrt{2c}x_n,\ldots,\sqrt{2c}x_1,c
angle$$

# Kernel Functions: Polynomial

In other words, using a quadratic kernel is equivalent to using a standard SVM where you've expanded the feature vectors to include:

- Each original feature value (times a constant)
- Each feature value squared
- The product of each pair of feature values (times a constant)
  - This can be especially useful, since it can capture *interactions* between features

Without the kernel trick, this large feature set would be computationally expensive to work with.

### **Kernel Functions**

In general, the kernel trick can create new features as nonlinear combinations of the old features

• Data that are not linearly separable in the original feature space might be separable in the new space



#### The radial basis function (RBF kernel) is: $k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma ||\mathbf{x}_i - \mathbf{x}_j||^2)$ squared Euclidean distance

- One of the most popular SVM kernels
- Related to the Gaussian/normal distribution
- Interpretation as expanded feature vector?
  - It actually maps to a feature vector with infinitely many features... so technically equivalent, but impossible to implement without using the kernel trick.



From: http://qingkaikong.blogspot.com/2016/12/machine-learning-8-support-vector.html

#### The radial basis function (RBF kernel) is: $k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma ||\mathbf{x}_i - \mathbf{x}_j||^2)$

- In addition to C,  $\gamma$  also affects overfitting
- Large  $\gamma \rightarrow$  small differences in distance between  $\boldsymbol{x_i}$  and  $\boldsymbol{x_i}$  are magnified
  - This will cause the classifier to fit the training data better, but may do worse on future data



From: http://qingkaikong.blogspot.com/2016/12/machine-learning-8-support-vector.html

# Kernel Methods: Summary (1)

- Kernel SVM is a reformulation of SVM that uses similarity between instances
  - To make a prediction for a new instance, need to calculate kernel function for the new instance and all training instances that are the support vectors
- Kernel SVM is equivalent to an SVM with a expanded feature set
  - Sometimes there is an intuitive interpretation of what the "new" features mean; sometimes not
- Kernel SVM with a linear kernel is equivalent to a standard SVM

# Kernel Methods: Summary (2)

- Kernels can be useful when your data has a small number of features and/or when the dataset is not linearly separable
- Some kernels are prone to overfitting
  - High degree polynomial; RBF with high scaling parameter
- Kernel SVM has additional hyperparameters you have to choose
  - Type of kernel
  - Parameters of kernel (e.g., d in polynomial,  $\gamma$  in RBF)

# Kernel Methods: Summary (3)

Also be aware that:

- Kernel methods not unique to SVM (invented long before, for perceptron), but popularized by it.
- Lots of other kernel functions not shown here, but these are the most common.
  - Specialized kernels exist for certain types of data (e.g., biological sequences, syntax trees)