# CS 1675: Intro to Machine Learning Support Vector Machines

Prof. Adriana Kovashka University of Pittsburgh October 23, 2018

# Plan for this lecture

- Linear Support Vector Machines
- Non-linear SVMs and the "kernel trick"
- Extensions and further details (briefly)
  - Soft-margin SVMs
  - Multi-class SVMs
  - Comparison: SVM vs logistic regression
- Why SVM solution is what it is (briefly)

#### Linear classifiers

• Find linear function to separate positive and negative examples



## Support vector machines



- Discriminative classifier based on optimal separating line (for 2d case)
- Maximize the margin between the positive and negative training examples

## Support vector machines

• Want line that maximizes the margin.



$$\mathbf{x}_i$$
 positive  $(y_i = 1)$ : $\mathbf{x}_i \cdot \mathbf{w} + b \ge 1$  $\mathbf{x}_i$  negative  $(y_i = -1)$ : $\mathbf{x}_i \cdot \mathbf{w} + b \le -1$ For support, vectors, $\mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$ 



Let 
$$\mathbf{w} = \begin{bmatrix} a \\ c \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$$

ax + cy + b = 0

#### Aside: Lines in R<sup>2</sup>



Let 
$$\mathbf{W} = \begin{bmatrix} a \\ c \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x \\ y \end{bmatrix}$$

ax + cy + b = 0  $\downarrow$   $\mathbf{w} \cdot \mathbf{x} + b = 0$ 

#### Aside: Lines in R<sup>2</sup>



Let 
$$\mathbf{w} = \begin{bmatrix} a \\ c \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$$

ax + cy + b = 0  $\downarrow$   $\mathbf{w} \cdot \mathbf{x} + b = 0$ 

#### Aside: Lines in R<sup>2</sup>



$$et \quad \mathbf{W} = \begin{vmatrix} a \\ c \end{vmatrix} \quad \mathbf{X} = \begin{vmatrix} x \\ y \end{vmatrix}$$

$$ax + cy + b = 0$$

$$\downarrow$$

$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

Kristen Grauman

## Support vector machines

Want line that maximizes the margin.



 $|\mathbf{x}_i \cdot \mathbf{w} + b|$ 

 $||\mathbf{w}||$ 

## Support vector machines

• Want line that maximizes the margin.



$$\mathbf{x}_i$$
 positive  $(y_i = 1)$ : $\mathbf{x}_i \cdot \mathbf{w} + b \ge 1$  $\mathbf{x}_i$  negative  $(y_i = -1)$ : $\mathbf{x}_i \cdot \mathbf{w} + b \le -1$ For support, vectors, $\mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$ Distance between point  
and line: $|\mathbf{x}_i \cdot \mathbf{w} + b|$   
 $||\mathbf{w}||$ Therefore, the margin is  $2 / ||\mathbf{w}||$ 

#### Finding the maximum margin line

- 1. Maximize margin  $2/||\mathbf{w}||$
- 2. Correctly classify all training data points:

 $\mathbf{x}_i$  positive  $(y_i = 1)$ :  $\mathbf{x}_i \cdot \mathbf{w} + b \ge 1$  $\mathbf{x}_i$  negative  $(y_i = -1)$ :  $\mathbf{x}_i \cdot \mathbf{w} + b \le -1$ 

Quadratic optimization problem:

Minimize 
$$\frac{1}{2} \mathbf{w}^T \mathbf{w}$$
  
Subject to  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \ge 1$  for  $\sum_{\text{train}}^{\text{One}} y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \ge 1$ 

One constraint for each training point.

Note sign trick.

#### Finding the maximum margin line



Finding the maximum margin line

- Solution:  $\mathbf{w} = \sum_{i} \alpha_{i} y_{i} \mathbf{x}_{i}$  MORE DETAILS NEXT TIME  $b = y_{i} - \mathbf{w} \cdot \mathbf{x}_{i}$  (for any support vector)
- Classification function:

$$f(x) = \operatorname{sign} (\mathbf{w} \cdot \mathbf{x} + \mathbf{b})$$
$$= \operatorname{sign} \left( \sum_{i} \alpha_{i} y_{i} \mathbf{x}_{i} \cdot \mathbf{x} + b \right)$$

If f(x) < 0, classify as negative, otherwise classify as positive.

- Notice that it relies on an *inner product* between the test point *x* and the support vectors *x<sub>i</sub>*
- (Solving the optimization problem also involves computing the inner products *x<sub>i</sub>* · *x<sub>j</sub>* between all pairs of training points)

#### Inner product

• The decision boundary for the SVM and its optimization depend on the inner product of two data points (vectors):

 $\mathbf{x}_{i}^{T}\mathbf{x}_{j}$ 

 $f(x) = \operatorname{sign} (\mathbf{w} \cdot \mathbf{x} + \mathbf{b})$  $= \operatorname{sign} \left( \sum_{i} \alpha_{i} y_{i} \mathbf{x}_{i} \cdot \mathbf{x} + \mathbf{b} \right)$ 

The inner product is equal

$$(\mathbf{x}_i^T \mathbf{x}) = \|\mathbf{x}_i\|^* \|\mathbf{x}_i\| \cos \theta$$

If the angle in between them is 0 then: If the angle between them is 90 then:  $(\mathbf{x}_i^T \mathbf{x}) = \|\mathbf{x}_i\|^* \|\mathbf{x}_i\|$  $(\mathbf{x}_i^T \mathbf{x}) = 0$ 

#### The inner product measures how similar the two vectors are

#### Example



## Solving for the alphas

- We know that for the support vectors, f(x) = 1 or -1 exactly
- Add a 1 in the feature representation for the bias
- The support vectors have coordinates and labels:
  - x1 = [0 1 1], y1 = -1
  - x2 = [-1 3 1], y2 = +1
  - x3 = [1 3 1], y3 = +1
- Thus we can form the following system of linear equations:

#### Solving for the alphas

• System of linear equations:

 $\alpha 1 \ y1 \ dot(x1, x1) + \alpha 2 \ y2 \ dot(x1, x2) + \alpha 3 \ y3 \ dot(x1, x3) = y1$  $\alpha 1 \ y1 \ dot(x2, x1) + \alpha 2 \ y2 \ dot(x2, x2) + \alpha 3 \ y3 \ dot(x2, x3) = y2$  $\alpha 1 \ y1 \ dot(x3, x1) + \alpha 2 \ y2 \ dot(x3, x2) + \alpha 3 \ y3 \ dot(x3, x3) = y3$ 

$$-2 * \alpha 1 + 4 * \alpha 2 + 4 * \alpha 3 = -1$$
  
-4 \* \alpha 1 + 11 \* \alpha 2 + 9 \* \alpha 3 = +1  
-4 \* \alpha 1 + 9 \* \alpha 2 + 11 \* \alpha 3 = +1

• Solution:  $\alpha 1 = 3.5$ ,  $\alpha 2 = 0.75$ ,  $\alpha 3 = 0.75$ 

## Solving for w, b; plotting boundary

We know w =  $\alpha_1 y_1 x_1 + ... + \alpha_N y_N x_N$  where N = # SVs Thus w = -3.5 \* [0 1 1] + 0.75 [-1 3 1] + 0.75 [1 3 1] = [0 1 -2]

Separating out weights and bias, we have:  $w = [0 \ 1]$  and b = -2

For SVMs, we used this eq for a line: ax + cy + b = 0 where w = [a c]

Thus 
$$ax + b = -cy \rightarrow y = (-a/c) x + (-b/c)$$

Thus y-intercept is -(-2)/1 = 2

The decision boundary is perpendicular to w and it has slope -0/1 = 0

#### Example



# Plan for this lecture

- Linear Support Vector Machines
- Non-linear SVMs and the "kernel trick"
- Extensions and further details (briefly)
  - Soft-margin SVMs
  - Multi-class SVMs
  - Comparison: SVM vs logistic regression
- Why SVM solution is what it is (briefly)

#### Nonlinear SVMs

• Datasets that are linearly separable work out great:



• But what if the dataset is just too hard?



• We can map it to a higher-dimensional space:



Andrew Moore

### Nonlinear SVMs

 General idea: the original input space can always be mapped to some higher-dimensional feature space where the training set is separable:



#### Nonlinear kernel: Example

• Consider the mapping  $\varphi(x) = (x, x^2)$ 



$$\varphi(x) \cdot \varphi(y) = (x, x^2) \cdot (y, y^2) = xy + x^2 y^2$$
$$K(x, y) = xy + x^2 y^2$$

#### The "kernel trick"

- The linear classifier relies on dot product between vectors  $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i \cdot \mathbf{x}_j$
- If every data point is mapped into highdimensional space via some transformation
   Φ: x<sub>i</sub> → φ(x<sub>i</sub>), the dot product becomes: K(x<sub>i</sub>, x<sub>j</sub>) = φ(x<sub>i</sub>) · φ(x<sub>j</sub>)
- A *kernel function* is similarity function that corresponds to an inner product in some expanded feature space
- The kernel trick: instead of explicitly computing the lifting transformation  $\varphi(\mathbf{x})$ , define a kernel function K such that:  $K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}_j)$

#### Examples of kernel functions

• Linear: 
$$K(x_i, x_j) = x_i^T x_j$$

Polynomials of degree up to d:

$$K(x_i, x_j) = (x_i^T x_j + 1)^d$$

Ш

112

Gaussian RBF:

$$K(x_i, x_j) = \exp(-\frac{\|x_i - x_j\|}{2\sigma^2})$$

Histogram intersection:

$$K(x_i, x_j) = \sum_k \min(x_i(k), x_j(k))$$

#### The benefit of the "kernel trick"

• Example: Polynomial kernel for 2-dim features

$$k(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^{\mathrm{T}} \mathbf{z})^{2} = (1 + x_{1}z_{1} + x_{2}z_{2})^{2}$$
  

$$= 1 + 2x_{1}z_{1} + 2x_{2}z_{2} + x_{1}^{2}z_{1}^{2} + 2x_{1}z_{1}x_{2}z_{2} + x_{2}^{2}z_{2}^{2}$$
  

$$= (1, \sqrt{2}x_{1}, \sqrt{2}x_{2}, x_{1}^{2}, \sqrt{2}x_{1}x_{2}, x_{2}^{2})(1, \sqrt{2}z_{1}, \sqrt{2}z_{2}, z_{1}^{2}, \sqrt{2}z_{1}z_{2}, z_{2}^{2})^{\mathrm{T}}$$
  

$$= \phi(\mathbf{x})^{\mathrm{T}}\phi(\mathbf{z}).$$
(7.42)

- ... lives in 6 dimensions
- With the kernel trick, we directly compute an inner product in 2-dim space, obtaining a scalar that we add 1 to and exponentiate

#### Is this function a kernel?

Problem:

- Checking if a given k : X × X → ℝ fulfills the conditions for a kernel is *difficult*:
- We need to prove or disprove

$$\sum_{i,j=1}^n t_i k(x_i, x_j) t_j \ge 0.$$

for any set  $x_1, \ldots, x_n \in \mathcal{X}$  and any  $t \in \mathbb{R}^n$  for any  $n \in \mathbb{N}$ . Workaround:

 It is easy to *construct* functions k that are positive definite kernels. 1) We can construct kernels from scratch:

- For any  $\varphi: \mathcal{X} \to \mathbb{R}^m$ ,  $k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathbb{R}^m}$  is a kernel.
- If  $d: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$  is a distance function, i.e.

• 
$$d(x, x') \ge 0$$
 for all  $x, x' \in \mathcal{X}$ ,  
•  $d(x, x') = 0$  only for  $x = x'$ ,  
•  $d(x, x') = d(x', x)$  for all  $x, x' \in \mathcal{X}$ ,  
•  $d(x, x') \le d(x, x'') + d(x'', x')$  for all  $x, x', x'' \in \mathcal{X}$ ,  
then  $k(x, x') := \exp(-d(x, x'))$  is a kernel.

2) We can construct kernels from other kernels:

- if k is a kernel and  $\alpha > 0$ , then  $\alpha k$  and  $k + \alpha$  are kernels.
- if  $k_1, k_2$  are kernels, then  $k_1 + k_2$  and  $k_1 \cdot k_2$  are kernels.

- 1. Select a kernel function.
- 2. Compute pairwise kernel values between labeled examples.
- Use this "kernel matrix" to solve for SVM support vectors & alpha weights.
- 4. To classify a new example: compute kernel values between new input and support vectors, apply alpha weights, check sign of output.

#### Example: Learning gender w/ SVMs



Moghaddam and Yang, Learning Gender with Support Faces, TPAMI 2002 Moghaddam and Yang, Face & Gesture 2000

#### Example: Learning gender w/ SVMs



Kristen Grauman

#### Example: Learning gender w/ SVMs



SVMs performed better than humans, at either resolution

ullet

Figure 6. SVM vs. Human performance

# Plan for this lecture

- Linear Support Vector Machines
- Non-linear SVMs and the "kernel trick"
- Extensions and further details (briefly)
  - Soft-margin SVMs
  - Multi-class SVMs
  - Comparison: SVM vs logistic regression
- Why SVM solution is what it is (briefly)



subject to 
$$y_i \boldsymbol{w}^T \boldsymbol{x}_i \geq 1$$
 ,  
 $\forall i = 1, \dots, N$ 

#### Soft-margin SVMs (allow misclassification)



subject to 
$$y_i \boldsymbol{w}^T \boldsymbol{x}_i \ge 1 - \xi_i,$$
  
 $\xi_i \ge 0, \quad \forall i = 1, \dots, N$ 

#### Slack variables in soft-margin SVMs



Figure from Bishop

#### Effect of margin size vs miscl. cost (c)

Training set



Misclassification ok, want large margin

Misclassification not ok

Image: Kent Munthe Caspersen

#### Effect of margin size vs miscl. cost (c)

Including test set A



Misclassification ok, want large margin

Misclassification not ok

Image: Kent Munthe Caspersen

#### Effect of margin size

Including test set B



Misclassification ok, want large margin

Misclassification not ok

Image: Kent Munthe Caspersen

#### Multi-class problems

# Instead of just two classes, we now have C classes

- E.g. predict which movie genre a viewer likes best
- Possible answers: action, drama, indie, thriller, etc.

#### Two approaches:

- One-vs-all
- One-vs-one

#### Multi-class problems

#### One-vs-all (a.k.a. one-vs-others)

- Train C classifiers
- In each, pos = data from class *i*, neg = data from classes other than *i*
- The class with the most confident prediction wins
- Example:
  - You have 4 classes, train 4 classifiers
  - 1 vs others: score 3.5
  - 2 vs others: score 6.2
  - 3 vs others: score 1.4
  - 4 vs other: score 5.5
  - Final prediction: class 2
- Issues?

#### Multi-class problems

#### One-vs-one (a.k.a. all-vs-all)

- Train C(C-1)/2 binary classifiers (all pairs of classes)
- They all vote for the label
- Example:
  - You have 4 classes, then train 6 classifiers
  - 1 vs 2, 1 vs 3, 1 vs 4, 2 vs 3, 2 vs 4, 3 vs 4
  - Votes: 1, 1, 4, 2, 4, 4
  - Final prediction is class 4

Let 
$$h(z) = max(0, 1 - z)$$

We have the objective to minimize  $\xi_i$  where:  $y_i \mathbf{w}^T \mathbf{x}_i \ge 1 - \xi_i$   $\xi_i \ge 1 - y_i \mathbf{w}^T \mathbf{x}_i$   $\xi_i \ge 0$  $\xi_i \ge max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)$ 

Then we can define a loss:

$$h(y_i \mathbf{w}^T \mathbf{x}_i) = max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)$$

and unconstrained SVM objective:

$$min_{\mathbf{w}}\frac{1}{2}\|\mathbf{w}\|^2 + C\sum_i^N h(y_i\mathbf{w}^T\mathbf{x}_i)$$

#### SVMs vs logistic regression

 When viewed from the point of view of regularized empirical loss minimization, SVM and logistic regression appear quite similar:

SVM: 
$$\sum_{i=1}^{n} \left( 1 - y_i \left[ w_0 + \mathbf{x}_i^T \mathbf{w}_1 \right] \right)^+ + \| \mathbf{w}_1 \|^2 / 2$$
  

$$\sum_{i=1}^{n} \underbrace{-\log P(y_i | \mathbf{x}, \mathbf{w})}_{-\log \sigma\left(y_i \left[ w_0 + \mathbf{x}_i^T \mathbf{w}_1 \right] \right)} + \| \mathbf{w}_1 \|^2 / 2$$

where  $\sigma(z) = (1 + \exp(-z))^{-1}$  is the logistic function.

(Note that we have transformed the problem maximizing the penalized log-likelihood into minimizing negative penalized log-likelihood.)

#### SVMs vs logistic regression

• The difference comes from how we penalize "errors":

Both: 
$$\sum_{i=1}^{n} \operatorname{Loss}\left(\underbrace{y_i \left[w_0 + \mathbf{x}_i^T \mathbf{w}_1\right]}_{i=1}\right) + \|\mathbf{w}_1\|^2/2$$



#### SVMs: Pros and cons

- Pros
  - Kernel-based framework is very powerful, flexible
  - Often a sparse set of support vectors compact at test time
  - Work very well in practice, even with very small training sample sizes
  - Solution can be formulated as a quadratic program (next time)
  - Many publicly available SVM packages: e.g. LIBSVM, LIBLINEAR, SVMLight
- Cons
  - Can be tricky to select best kernel function for a problem
  - Computation, memory
    - At training time, must compute kernel values for all example pairs
    - Learning can take a very long time for large-scale problems