

# Unit 7

## Maximum-Margin Classification and the Support Vector Machine

Prof. Phil Schniter



THE OHIO STATE UNIVERSITY

ECE 4300: Introduction to Machine Learning, Sp20

# Learning objectives

- Gain experience with linear classification of images
  - representing images, displaying images, classifying images
- Understand the geometry of the linear classification boundary:
  - orthogonality to  $w$ , offset  $b$ , margin  $1/\|w\|$
- Understand the margin-maximizing classifier
- Understand the support vector classifier (SVC)
- Understand the support vector machine (SVM)
- Understand how to implement the SVC and SVM with `sklearn`

# Outline

- Motivating Example: Recognizing Handwritten Digits
- Maximum-Margin Classification
- The Support Vector Classifier
- The Support Vector Machine

# MNIST digit classification

- Goal: Handwriting recognition
- Dataset: **MNIST**
  - Only digits 0–9
  - 28x28 images (grayscale)
  - 70,000 examples
  - Collected by American Census Bureau in 1980s
- Why use it?
  - Simple
  - Widely accessible
  - Many benchmarks



```
from sklearn.datasets import fetch_openml
mnist = fetch_openml('mnist_784', version=1, cache=True)
mnist.data.shape
```

```
(70000, 784)
```

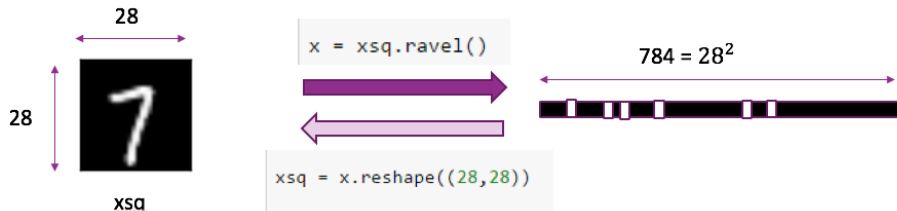
# MNIST benchmarks

- In this unit, we focus on the **support vector machine** (SVM).
- Not state-of-the-art, but quite good, and widely used in machine learning

Type ↕	Classifier ↕	Distortion ↕	Preprocessing ↕	Error rate (%) ▼
Linear classifier	Pairwise linear classifier	None	Deskewing	7.6 <sup>[9]</sup>
Non-linear classifier	40 PCA + quadratic classifier	None	None	3.3 <sup>[9]</sup>
Deep <b>Neural network</b>	2-layer 784-800-10	None	None	1.6 <sup>[21]</sup>
Boosted Stumps	Product of stumps on <b>Haar features</b>	None	Haar features	0.87 <sup>[19]</sup>
Deep <b>Neural network</b>	2-layer 784-800-10	elastic distortions	None	0.7 <sup>[21]</sup>
<b>Support vector machine</b>	Virtual <b>SVM</b> , deg-9 poly, 2-pixel jittered	None	Deskewing	0.56 <sup>[20]</sup>
<b>K-Nearest Neighbors</b>	K-NN with non-linear deformation (P2DHMDM)	None	Shiftable edges	0.52 <sup>[18]</sup>
Deep <b>neural network</b>	6-layer 784-2500-2000-1500-1000-500-10	elastic distortions	None	0.35 <sup>[22]</sup>
<b>Convolutional neural network</b>	6-layer 784-40-80-500-1000-2000-10	None	Expansion of the training data	0.31 <sup>[15]</sup>
<b>Convolutional neural network</b>	6-layer 784-50-100-500-1000-10-10	None	Expansion of the training data	0.27 <sup>[23]</sup>
<b>Convolutional neural network</b>	Committee of 35 CNNs, 1-20-P-40-P-150-10	elastic distortions	Width normalizations	0.23 <sup>[8]</sup>
<b>Convolutional neural network</b>	Committee of 5 CNNs, 6-layer 784-50-100-500-1000-10-10	None	Expansion of the training data	0.21 <sup>[17]</sup>

# Representing images in Numpy

- Each pixel is a value between 0 (black) and 255 (white)
- Images can be represented as matrices or vectors



- In **Numpy**, vectorization is done by stacking rows, but in most papers/books it is done by stacking columns:

$$\mathbf{X} = \begin{bmatrix} x_1 & x_{29} & \cdots & x_{757} \\ x_2 & x_{30} & \cdots & x_{758} \\ \vdots & \vdots & & \vdots \\ x_{28} & x_{56} & \cdots & x_{784} \end{bmatrix} = \text{mat}(\mathbf{x}), \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{784} \end{bmatrix} = \text{vec}(\mathbf{X})$$

- Basically, **Numpy** works with  $\mathbf{X}^T$  and  $\mathbf{x}^T$  for the  $\mathbf{X}$  and  $\mathbf{x}$  defined above

# Displaying images with matplotlib

- `plt.imshow` is the main plot command
- We randomly permute the sample indices using `Iperm`
  - Else training set might include different digits than testing set
  - Do this for any dataset!
- A random sample of 4 MNIST digits:



- A human could easily classify these

```
def plt_digit(x):
    nrow = 28
    ncol = 28
    xsq = x.reshape((nrow,ncol))
    plt.imshow(xsq, cmap='Greys_r')
    plt.xticks([])
    plt.yticks([])

# Convert data to a matrix
X = mnist.data
y = mnist.target.astype(np.int8) # fetch_op

# Select random digits
nplt = 4
nsamp = X.shape[0]
Iperm = np.random.permutation(nsamp)

# Plot the images using the subplot command
for i in range(nplt):
    ind = Iperm[i]
    plt.subplot(1,nplt,i+1)
    plt_digit(X[ind,:])
```

# Try multinomial logistic regression (MLR)

- Train and test on 10000 samples
  - More training leads to better results
  - But can be very slow
- Select a fast solver (**lbfgs**)
  - Even this can take minutes

```
ntr = 10000
nts = 10000
Xs = X/255.0*2 - 1 # convert to range [-1,1]
Xtr = Xs[Iperm[:ntr],:]
ytr = y[Iperm[:ntr]]
Xts = Xs[Iperm[ntr:ntr+nts],:]
yts = y[Iperm[ntr:ntr+nts]]
```

```
logreg = linear_model.LogisticRegression(verbose=1, multi_class='multinomial',
                                          solver='lbfgs', max_iter=500)
```

```
logreg.fit(Xtr,ytr)
```

```
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers
/Users/schniter.1/anaconda3/lib/python3.6/site-packages/sklearn/linear_model/logit.py:469:
ConvergenceWarning: lbfgs failed to converge. Increase the number of iterations
"of iterations.", ConvergenceWarning)
[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 7.2s finished
```

```
LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
                   intercept_scaling=1, max_iter=500, multi_class='multinomial',
                   n_jobs=None, penalty='l2', random_state=None, solver='lbfgs',
                   tol=0.0001, verbose=1, warm_start=False)
```



# Performance of MLR

- Accuracy = 89% (not very good)
- Most of the errors MLR made would be obvious to a human, for example:

```
yhat = logreg.predict(Xts)  
acc = np.mean(yhat == yts)  
print('Accuracy = {0:f}'.format(acc))
```

Accuracy = 0.892900

true=5 est=6



true=0 est=5



true=5 est=3



true=3 est=2

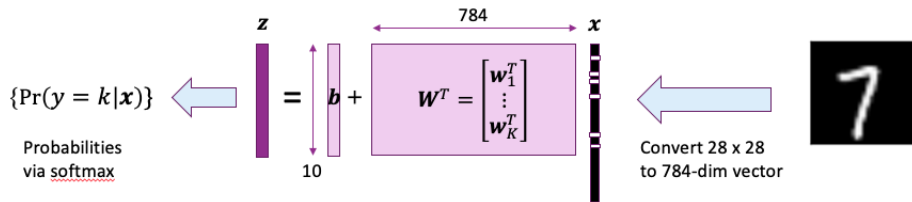


- What went wrong?
- Can we do better?

# Review of MLR

- To classify a test vector  $x$  (e.g., an MNIST image), we do the following:
- For each class hypothesis  $k = 1, \dots, K$ ,
  - design weights  $b_k$  and  $w_k = [w_{k1}, \dots, w_{kd}]^T$  that classify  $x$  as **in- $k$**  or **not-in- $k$**
  - compute a linear **score**  $z_k = b_k + x^T w_k$

then classify as  $\hat{y} = \arg \max_k \underbrace{\Pr\{y=k | x\}}_{\text{probability of } x \text{ in class } k} = \arg \max_k \underbrace{\frac{e^{z_k}}{\sum_{l=1}^K e^{z_l}}}_{\text{softmax}} = \arg \max_k z_k$



# What weights does MLR learn?

- Intuitively, we want  $z_k = b_k + \mathbf{x}^\top \mathbf{w}_k$  to be large when  $\mathbf{x}$  is from class  $k$ , and small otherwise.
- Let's visualize the learned  $\mathbf{w}_k$  for  $k = 0, \dots, 4$  by plotting them as images:

Learned from  
100 training  
samples:



Learned from  
10000 training  
samples:



- MLR struggles to match *all* the training examples from each class due to variations across that class, especially when number of training samples is large!

# Outline

- Motivating Example: Recognizing Handwritten Digits
- **Maximum-Margin Classification**
- The Support Vector Classifier
- The Support Vector Machine

# Hyperplane geometry

- In *binary* linear classification, the decision boundary is  $z = 0$  or

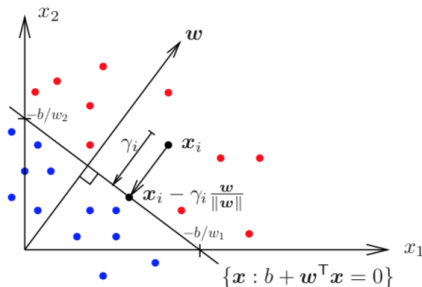
$$\{x : b + w^T x = 0\}$$

which is a **hyperplane**

- This **hyperplane is orthogonal to  $w$** :

- For any  $\{x_1, x_2\}$  on the hyperplane,  
 $b + w^T x_i = 0$  for  $i = 1, 2$ .

- Thus  $w^T(x_1 - x_2) = 0$



- From the figure, the (signed) distance from a point  $x_i$  to the hyperplane is  $\gamma_i$ :

$$0 = b + \left(x_i - \gamma_i \frac{w}{\|w\|}\right)^T w = b + x_i^T w - \gamma_i \|w\| \quad \Rightarrow \quad \boxed{\gamma_i = \frac{b + w^T x_i}{\|w\|}}$$

- The classifier  $(\alpha b, \alpha w)$  is **equivalent** to the classifier  $(b, w)$  for all  $\alpha > 0$

# Review of linear separability

- Assume binary training data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$
- The data is “**linearly separable**” if there exists  $(b, \mathbf{w})$  such that

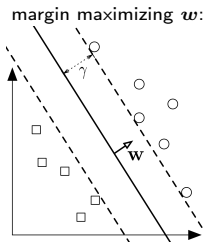
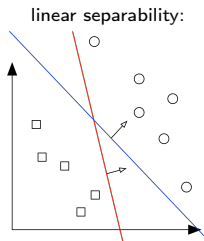
$$b + \mathbf{x}_i^T \mathbf{w} > 1 \quad \text{whenever } y_i = 1$$

$$b + \mathbf{x}_i^T \mathbf{w} < -1 \quad \text{whenever } y_i = -1$$

or, in a more compact form,

$$y_i(b + \mathbf{x}_i^T \mathbf{w}) > 1 \quad \forall i$$

- Note:  $(b, \mathbf{w})$  is not unique!
- “1” chosen for convenience; could use any value
- If separable, the distance  $\gamma_i$  from the boundary to datapoint  $\mathbf{x}_i$  obeys  $y_i \gamma_i \geq 1/\|\mathbf{w}\| \triangleq \gamma$ , where  $\gamma$  is known as the “**margin**”



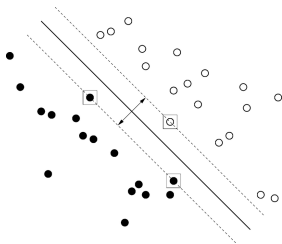
# Maximum-margin classification

- To make the classifier robust, we should *maximize the margin*  $1/\|w\|$ !
  - Equivalently, we should minimize  $\|w\|$

- The **hard-margin classifier** is defined as the  $(b, w)$  that solves

$$\min \frac{1}{2} \|w\|^2 \quad \text{s.t.} \quad y_i(b + x_i^\top w) \geq 1 \quad \forall i = 1, \dots, n$$

- The orientation of its boundary depends only on a few samples  $x_i$  near the boundary:
  - Called the **support vectors**
  - More on this later ...



- Problem: This classifier exists only if the data is linear separable!

# Outline

- Motivating Example: Recognizing Handwritten Digits
- Maximum-Margin Classification
- The Support Vector Classifier
- The Support Vector Machine



# The support vector classifier

- To handle data that is not linearly separable, we replace the hard constraint

$$y_i(b + \mathbf{x}_i^T \mathbf{w}) \geq 1 \quad \forall i \quad \Leftrightarrow \quad 0 \geq 1 - y_i(b + \mathbf{x}_i^T \mathbf{w}) \quad \forall i$$

with the cost function 
$$\underbrace{\sum_{i=1}^n \max \{0, 1 - y_i(b + \mathbf{x}_i^T \mathbf{w})\}}_{= 0 \text{ iff linearly separable}}$$

- This results in the **soft-margin classifier** or **support-vector classifier (SVC)**:

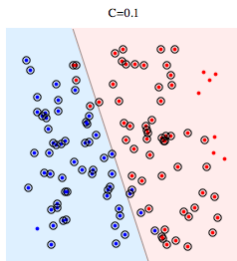
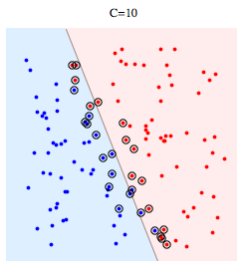
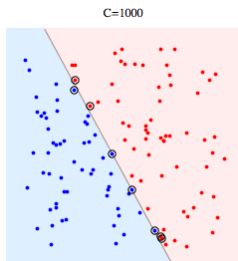
$$\min_{b, \mathbf{w}} \left\{ C \sum_{i=1}^n \max \{0, 1 - y_i(b + \mathbf{x}_i^T \mathbf{w})\} + \frac{1}{2} \|\mathbf{w}\|^2 \right\} \text{ with tunable parameter } C > 0$$

- $C$  trades off between margin size  $1/\|\mathbf{w}\|$  and number of margin violations
  - Larger  $C$  penalizes violations more and thus yields a smaller margin  $1/\|\mathbf{w}\|$
  - For linearly separable data and sufficiently large  $C$ , soft-margin  $\Leftrightarrow$  hard-margin

# Effect of SVC penalty $C$

$C$  controls the margin  $\frac{1}{\|w\|}$ , and thus the number of support vectors

- Larger  $C$ : smaller margin, fewer support vectors
  - More sensitive to datapoints close to hyperplane: Lower bias, higher variance
- Smaller  $C$ : larger margin, more support vectors
  - Less sensitive to data: Higher bias, lower variance



Support  
vectors  
are circled

The support vectors are the  $x_i$  with signed distance  $< \frac{1}{\|w\|}$  from hyperplane

# The SVC and hinge loss

- Recalling L2-regularized logistic regression:

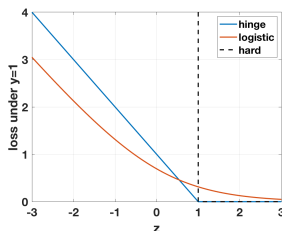
$$\min_{\mathbf{w}, b} \left\{ \sum_{i=1}^n \underbrace{(\ln[1 + e^{z_i}] - y_i z_i)}_{\text{"logistic loss"}} + \frac{1}{2C} \|\mathbf{w}\|^2 \right\} \quad \text{for } z_i = b + \mathbf{x}_i^T \mathbf{w}$$

we see that the SVC looks similar, but with a different loss:

$$\min_{\mathbf{w}, b} \left\{ \sum_{i=1}^n \underbrace{\max\{0, 1 - y_i z_i\}}_{\text{"hinge loss"}} + \frac{1}{2C} \|\mathbf{w}\|^2 \right\} \quad \text{for } z_i = b + \mathbf{x}_i^T \mathbf{w}$$

- Comparison of loss functions:

- hinge loss places *no penalty* on samples lying on good side of margin (i.e.,  $y_i z_i > 1$ )
- thus, the SVC uses only the samples  $\mathbf{x}_i$  that cause  $y_i z_i \leq 1$ , called "**support vectors**"
- The SVC ignores the other samples!



# A closer look at the SVC

- We just saw that the SVC solution  $(\mathbf{w}_*, b_*)$  minimizes the cost

$$J_{\text{svc}}(\mathbf{w}, b) \triangleq C \sum_{i=1}^n \max \{0, 1 - y_i \underbrace{(b + \mathbf{w}^\top \mathbf{x}_i)}_{z_i}\} + \frac{1}{2} \|\mathbf{w}\|^2$$

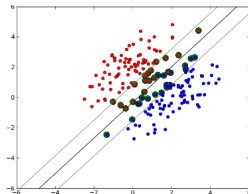
- The gradient must equal zero at  $(\mathbf{w}_*, b_*)$ :

$$0 = \frac{\partial}{\partial b} J(\mathbf{w}_*, b_*) = -C \sum_{i=1}^n \alpha_i y_i \quad \text{for } \alpha_i = \begin{cases} 1 & y_i(b_* + \mathbf{w}_*^\top \mathbf{x}_i) < 1 \\ ? & y_i(b_* + \mathbf{w}_*^\top \mathbf{x}_i) = 1 \\ 0 & y_i(b_* + \mathbf{w}_*^\top \mathbf{x}_i) > 1 \end{cases}$$

$$\mathbf{0} = \nabla_{\mathbf{w}} J(\mathbf{w}_*, b_*) = -C \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i + \mathbf{w}_*$$

which implies  $\mathbf{w}_* = C \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$  and  $\sum_{i=1}^n \alpha_i y_i = 0$

- The equation for  $\mathbf{w}_*$  shows that it depends only on the **support vectors**  $\{\mathbf{x}_i : y_i(b_* + \mathbf{x}_i^\top \mathbf{w}_*) \leq 1\}$

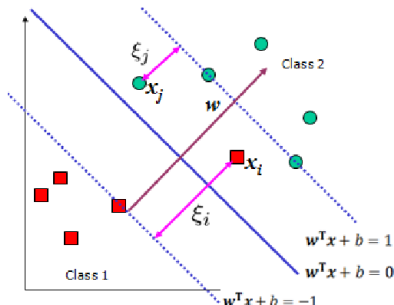


# Common form of the SVC

- The SVC is often written using “slack variables”  $\xi_i \triangleq \max\{0, 1 - y_i z_i\}$ :

$$\min_{b, \mathbf{w}} \left\{ C \sum_{i=1}^n \xi_i + \frac{1}{2} \|\mathbf{w}\|^2 \right\} \text{ s.t. } y_i(b + \mathbf{x}_i^T \mathbf{w}) \geq 1 - \xi_i \text{ and } \xi_i \geq 0 \text{ for all } i$$

- $\xi_i$  measures violation of the margin
  - $\xi_i = 0$ : sample outside “street” (on correct side of hyperplane)
  - $\xi_i \in (0, 2)$ : sample in “street”
  - $\xi_i > 1$ : sample misclassified (wrong side of hyperplane)
- Any  $\xi_i > 0$  is considered a violation



# Applying the SVC to test data

- Given a test vector  $\mathbf{x}$ , the SVC computes the score

$$z = b_* + \mathbf{w}_*^T \mathbf{x} \quad \text{with} \quad \mathbf{w}_* = C \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \quad \text{and} \quad \alpha_i = 0 \text{ for non-support } \mathbf{x}_i$$

- The resulting score

$$z = b_* + C \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^T \mathbf{x} \quad \text{used for} \quad \hat{y} = \text{sgn}(z)$$

is a weighted linear combination of the “correlations”  $\mathbf{x}_i^T \mathbf{x}$  between  $\mathbf{x}$  and each support vector  $\mathbf{x}_i$ .

- The SVC uses only  $\{\mathbf{x}_i\}$  that are support vectors, whereas logistic regression uses *all* the training samples  $\{\mathbf{x}_i\}$
- The SVC tries to **focus on what is most important** when making a decision

# Multiclass SVC

- In multiclass linear classification with  $K$  classes, we compute the score vector

$$\mathbf{z}_i = \begin{bmatrix} b_1 \\ \vdots \\ b_K \end{bmatrix} + \begin{bmatrix} \mathbf{w}_1^\top \\ \vdots \\ \mathbf{w}_K^\top \end{bmatrix} \mathbf{x}_i = \mathbf{b} + \mathbf{W}^\top \mathbf{x}_i \in \mathbb{R}^K$$

- One option is to *separately* design one-vs-all SVCs  $\{(b_k, \mathbf{w}_k)\}_{k=1}^K$

- Another option** is to *jointly* design the multiclass SVC  $(\mathbf{b}, \mathbf{W})$  via

$$\min_{\mathbf{b}, \mathbf{W}} \left\{ \frac{1}{n} \sum_{i=1}^n \sum_{j \neq y_i} \max \{0, 1 + z_{ij} - z_{i, y_i}\} + \frac{1}{2C} \|\mathbf{W}\|_F^2 \right\}$$

- Note that, when  $K = 2$ , we get (using  $j, y_i \in \{-1, 1\}$ )

$$\begin{aligned} \sum_{j \neq y_i} \max \{0, 1 + z_{ij} - z_{i, y_i}\} &= \max \{0, 1 + z_{i, -y_i} - z_{i, y_i}\} \\ &= \max \{0, 1 - y_i(z_{i, 1} - z_{i, -1})\} \\ &= \max \{0, 1 - y_i(\mathbf{b} + \mathbf{w}^\top \mathbf{x}_i)\} \end{aligned}$$

for  $b \triangleq b_1 - b_{-1}$  and  $\mathbf{w} \triangleq \mathbf{w}_1 - \mathbf{w}_{-1}$ , which matches our earlier binary SVC

# Implementing the SVC in sklearn

- The SVC is easy to implement in `sklearn` using `svm.SVC`
  - The parameter  $C$  penalizes the slack variables
  - As we discussed, larger  $C$  means fewer support vectors

```
from sklearn import svm
svc = svm.SVC(probability=False, kernel="linear", verbose=1)
svc.fit(Xtr,ytr)
yhat_ts = svc.predict(Xts)
acc = np.mean(yhat_ts == yts)
print('Accuracy = {0:f}'.format(acc))
```

Accuracy = 0.907100

- The SVC outperforms logistic regression by about 2% on this dataset
  - good, but not a huge improvement
- For large datasets, it is recommended to use `svm.LinearSVC` instead

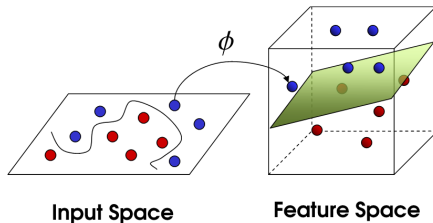


# Outline

- Motivating Example: Recognizing Handwritten Digits
- Maximum-Margin Classification
- The Support Vector Classifier
- The Support Vector Machine

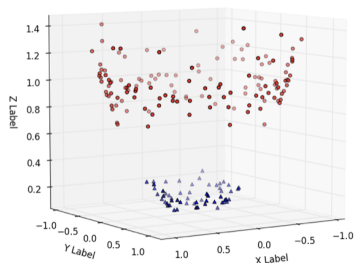
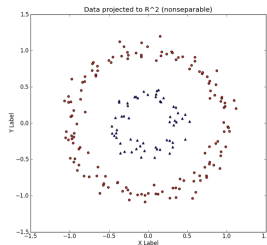
# Support vector machine

- Support vector classifier (SVC)
  - works very well when the features are roughly linearly separable
- Support vector machine (SVM)
  - SVC plus ...
    - 1) a **feature transformation**  $x \mapsto \phi(x)$  to enhance linear separability
    - 2) the “**kernel trick**” for fast implementation



# Feature transformations

- We want a feature transformation  $x \mapsto \phi(x)$  that leads to linear separability
- Usually we do this by mapping to a higher-dimensional space
- Example: For the  $x = [x_1, x_2]^T$  data below, which is not linearly separable, we can use  $\phi(x) = [x_1, x_2, x_1^2 + x_2^2]^T$  to get linear separability!



# The problem of dimensionality

- In principle, given transformed features  $\phi(x)$ , we could train via

$$(\mathbf{w}_*, b_*) = \min_{\mathbf{w}, b} \left\{ C \sum_{i=1}^n \max \{0, 1 - y_i z_i\} + \frac{1}{2} \|\mathbf{w}\|^2 \right\} \quad \text{for } z_i = b + \mathbf{w}^\top \phi(\mathbf{x}_i)$$

and classify a test vector  $x$  via

$$\hat{y} = \text{sgn}(z) = \text{sgn} \left[ b_* + C \sum_{i=1}^n \alpha_i y_i \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}) \right]$$

- But what if the dimension of  $\phi(x)$  is very high?
  - The dimensionality of  $\mathbf{w}$  will be equally high
  - The optimization problem will be computationally expensive!
- And what if  $\phi(x)$  is infinite-dimensional?
  - The optimization problem can't be solved using the above formulation
  - This might sound crazy, but it's not unusual (more details soon ...)

# The Lagrangian dual formulation of the SVC

- Can avoid these problems by using the “**Lagrangian dual**” form of the SVC:

$$\mathbf{w}_* = \sum_{i=1}^n \lambda_i y_i \mathbf{x}_i \quad \text{and} \quad b_* = \mathbf{w}_*^T \mathbf{x}_{i_*} - y_{i_*} \quad \text{for any } i_* \text{ s.t. } \lambda_{i_*} > 0$$

where the Lagrange multipliers  $\{\lambda_i\}$  solve the quadratic programming problem

$$\min_{\lambda} \left\{ \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \right\} \text{ s.t. } \lambda_i \in [0, C] \text{ and } \sum_{i=1}^n \lambda_i y_i = 0$$

- The derivation is outside the scope of this course

See Sec. 12.2 of **Hastie, Tibshirani, Friedman, *The Elements of Statistical Learning, 2nd Ed.*, 2009**

- But some terms may look familiar from our gradient analysis of the SVC cost
- Note that we can avoid computing  $\mathbf{w}_*$  by directly computing the decision

$$\hat{y} = \text{sgn} \left( b_* + \sum_{i=1}^n \lambda_i y_i \mathbf{x}_i^T \mathbf{x} \right) \quad \text{where } b_* = \sum_{i=1}^n \lambda_i y_i \mathbf{x}_i^T \mathbf{x}_{i_*} - y_{i_*}$$

# The kernel trick

- Now consider the dual-form SVC with  $\mathbf{x}$  replaced by the transformation  $\phi(\mathbf{x})$ :

$$\hat{y} = \text{sgn} \left( b_* + \sum_{i=1}^n \lambda_i y_i \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}) \right) \text{ where } b_* = \sum_{i=1}^n \lambda_i y_i \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_{i_*}) - y_{i_*}$$

where  $i_*$  is any support index, and where  $\{\lambda_i\}$  solve

$$\min_{\lambda} \left\{ \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) \right\} \text{ s.t. } \begin{cases} \lambda_i \in [0, C] \\ \sum_{i=1}^n \lambda_i y_i = 0 \end{cases}$$

- The transformed features appear only in the form of a “kernel” function

$$\kappa(\mathbf{x}, \mathbf{x}') \triangleq \phi(\mathbf{x})^\top \phi(\mathbf{x}')$$

- Often the kernel is easier to compute than  $\phi(\mathbf{x})$
- For both learning and inference, we only need to compute the kernel (not  $\phi(\mathbf{x})$ )!
- Called the “kernel trick”

# Popular kernels

- Popular kernels include:
  - **radial basis function (RBF)**:  $\kappa(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$  for “width”  $\gamma > 0$
  - **polynomial**:  $\kappa(\mathbf{x}, \mathbf{x}') = (\gamma \mathbf{x}^\top \mathbf{x}' + r)^d$  where  $\gamma > 0$ ,  $r \geq 0$ , and often  $d = 2$
  - **sigmoidal**:  $\kappa(\mathbf{x}, \mathbf{x}') = \tanh(\gamma \mathbf{x}^\top \mathbf{x}' + r)$  for some  $\gamma > 0$  and  $r > 0$
  - **linear**:  $\kappa(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}' \dots$  corresponds to trivial transformation  $\phi(\mathbf{x}) = \mathbf{x}$
- Each kernel measures a “similarity” between  $\mathbf{x}$  and  $\mathbf{x}'$
- Note: The feature space of the RBF kernel is infinite dimensional!

$$\exp(-\tfrac{1}{2}\|\mathbf{x} - \mathbf{x}'\|^2) = \exp(-\tfrac{1}{2}\|\mathbf{x}\|^2) \exp(-\tfrac{1}{2}\|\mathbf{x}'\|^2) \sum_{j=0}^{\infty} \frac{(\mathbf{x}^\top \mathbf{x}')^j}{j!}$$

- Can't implement the RBF kernel via  $\phi(\mathbf{x})$
- But can implement RBF efficiently using the kernel trick!

# Implementing the SVM in sklearn

- The SVM is easy to implement in **sklearn** using **svm.SVC**
  - Choose from several kernels: **rbf** (default), **poly**, **sigmoid**, **linear**, or custom.
  - Specify the kernel width  $\gamma > 0$
  - Control number of support vectors via  $C > 0$
  - Ideally optimize all parameters via cross-validation
    - See <http://peekaboo-vision.blogspot.com/2010/09/mnist-for-ever.html>
    - Below we use settings from <https://martin-thoma.com/svm-with-sklearn/>

```
from sklearn import svm
svcrbf = svm.SVC(probability=False, kernel="rbf", C=2.8, gamma=.0073, verbose=1)
svcrbf.fit(Xtr,ytr)
yhat_ts = svcrbf.predict(Xts)
acc = np.mean(yhat_ts == yts)
print('Accuracy = {0:f}'.format(acc))
```

Accuracy = 0.970500

- The SVM significantly outperforms both logistic regression and SVC!
  - Could do even better with more training samples and tweaked parameters



# Learning objectives

- Gain experience with linear classification of images
  - representing images, displaying images, classifying images
- Understand the geometry of the linear classification boundary:
  - orthogonality to  $w$ , offset  $b$ , margin  $1/\|w\|$
- Understand the margin-maximizing classifier
- Understand the support vector classifier (SVC)
- Understand the support vector machine (SVM)
- Understand how to implement the SVC and SVM with `sklearn`