

CS5489 - Machine Learning

Lecture 6c - Non-Linear Dimensionality Reduction

Dr. Antoni B. Chan

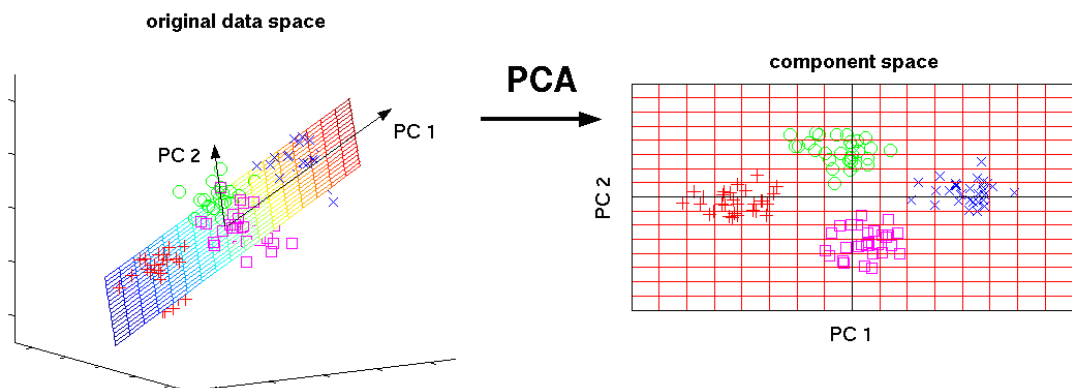
Dept. of Computer Science, City University of Hong Kong

Outline

1. Linear Dimensionality Reduction for Vectors
2. Linear Dimensionality Reduction for Text
3. **Non-linear Dimensionality Reduction**
4. Manifold Embedding

Linear Dimensionality Reduction

- PCA, NMF, LSA are all linear dimensionality reduction methods
 - model the data as "living" on a linear manifold (line, plane, etc).

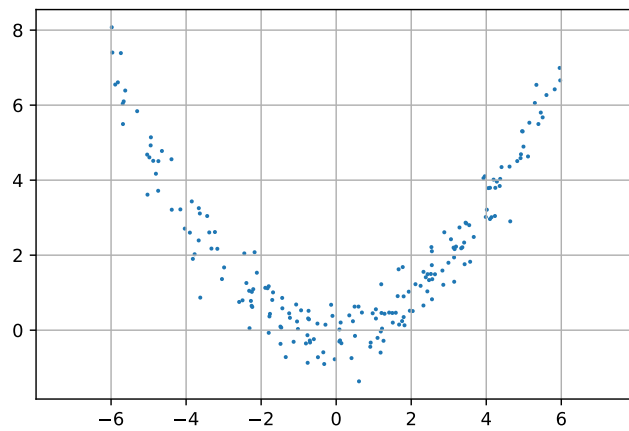


Non-linear surface

- What if the data "lives" on a non-flat surface?

```
In [4]: pfig
```

```
Out[4]:
```

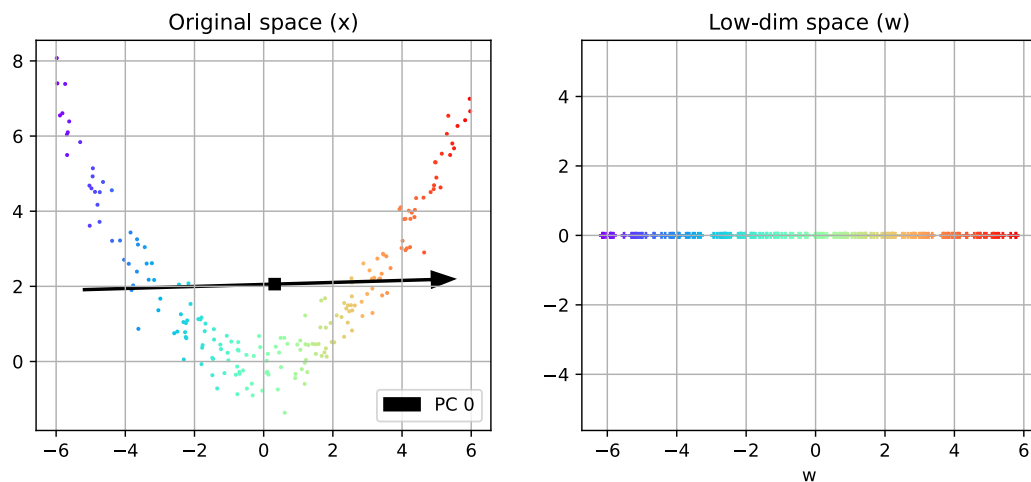


- PCA can't capture the curvature of the data
 - purple points are close together
 - red points are close together

In [5]:

```
pca = decomposition.PCA(n_components=1)
W = pca.fit_transform(X)

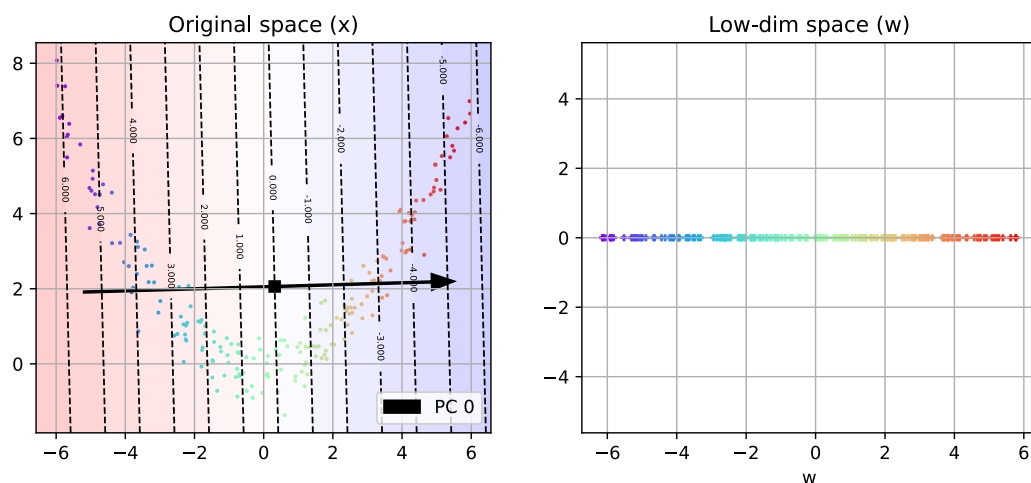
plt.figure(figsize=(10,4))
plot_basis(X, pca.components_, Y=Y, showlowerarrow=False)
```



- iso-contours of PCA projection
 - points on the same dashed line are projected to the same PCA coefficient.

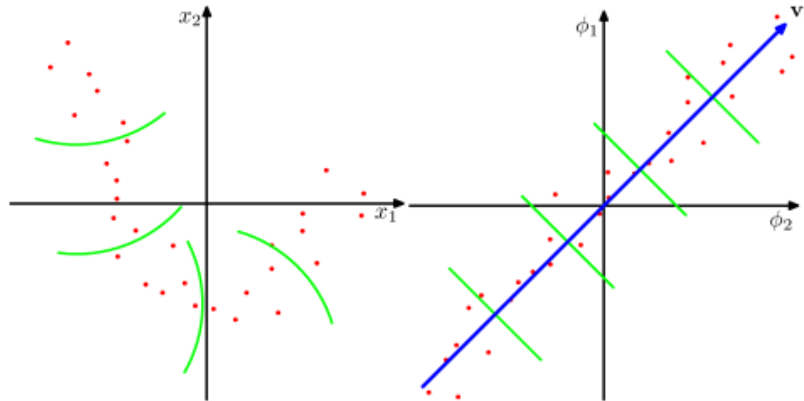
In [6]:

```
plt.figure(figsize=(10,4))
plot_basis(X, pca.components_, Y=Y, showcontours=True, pca=pca, showlowerarrow=False)
```



Kernel PCA

- How to project to a non-linear surface?
 - apply a high-dimensional feature transformation to the data
 - $\mathbf{x}_i \Rightarrow \phi(\mathbf{x}_i)$
 - project high-dim data to a linear surface
 - i.e. run PCA on $\phi(\mathbf{x}_i)$
 - in the original space, the projection will be non-linear



Kernel principal components

- a principal component \mathbf{v} is a linear combination of high-dim vectors
 - $\mathbf{v} = \sum_{i=1}^n a_i \phi(\mathbf{x}_i)$
 - where a_i are learned weights for component \mathbf{v} .
- For a new point \mathbf{x}_* , the KPCA coefficient for \mathbf{v} is
 - $w = \phi(\mathbf{x}_*)^T \mathbf{v} = \sum_{i=1}^n a_i \phi(\mathbf{x}_*)^T \phi(\mathbf{x}_i) = \sum_{i=1}^n a_i k(\mathbf{x}_*, \mathbf{x}_i) = \mathbf{k}_*^T \mathbf{a}$
 - coefficient is based on similarity to data points belonging to \mathbf{v} .
 - using the kernel trick saves computation.

Learning KPCA weights

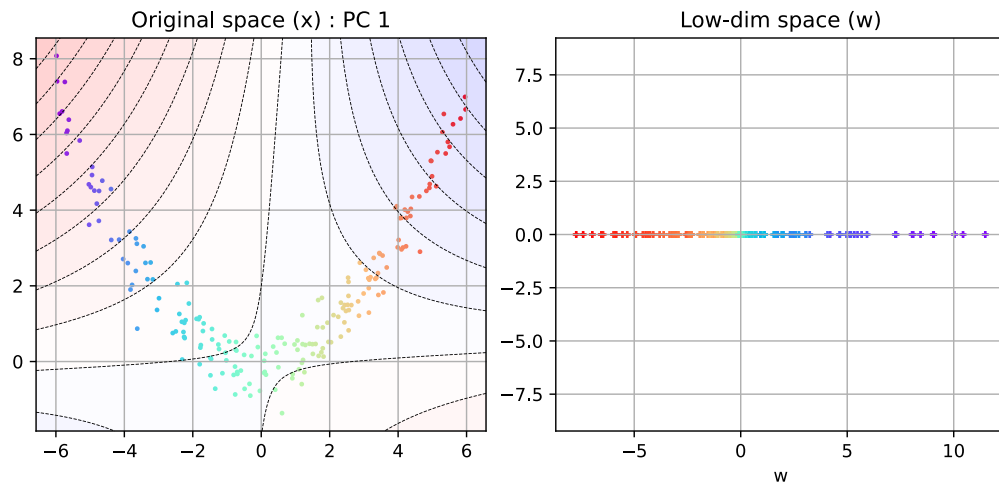
- Apply PCA to the transformed data $\{\phi(\mathbf{x}_i)\}$:
- After some derivation:
 - 1) Calculate the kernel matrix $\mathbf{K} = [k(\mathbf{x}_i, \mathbf{x}_j)]_{ij}$
 - 2) center the kernel (remove the mean in the high-dim space):

$$\tilde{\mathbf{K}} = (\mathbf{I} - \frac{1}{N} \mathbf{1}\mathbf{1}^T) \mathbf{K} (\mathbf{I} - \frac{1}{N} \mathbf{1}\mathbf{1}^T)$$
 - $\mathbf{1}\mathbf{1}^T$ is a matrix of ones.
 - 3) Find the top-K eigenvector/value pairs: $\tilde{\mathbf{K}} \mathbf{a}_j = \lambda_j \mathbf{a}_j$
 - 4) Scale (normalize PC in high-dim space): $\mathbf{a}_j \leftarrow \frac{1}{\sqrt{\lambda_j}} \mathbf{a}_j$
 - 5) Project new data \mathbf{x}_* : $w_j = \mathbf{k}_*^T \mathbf{a}_j$
- Example using polynomial kernel
 - purple points are further apart.
 - PC coefficient corresponds to location along the data curve.

In [8]:

```
# run KPCA
kpca = decomposition.KernelPCA(n_components=1, kernel='poly', gamma=0.15, degree=2, coef0=0)
W = kpca.fit_transform(X)
```

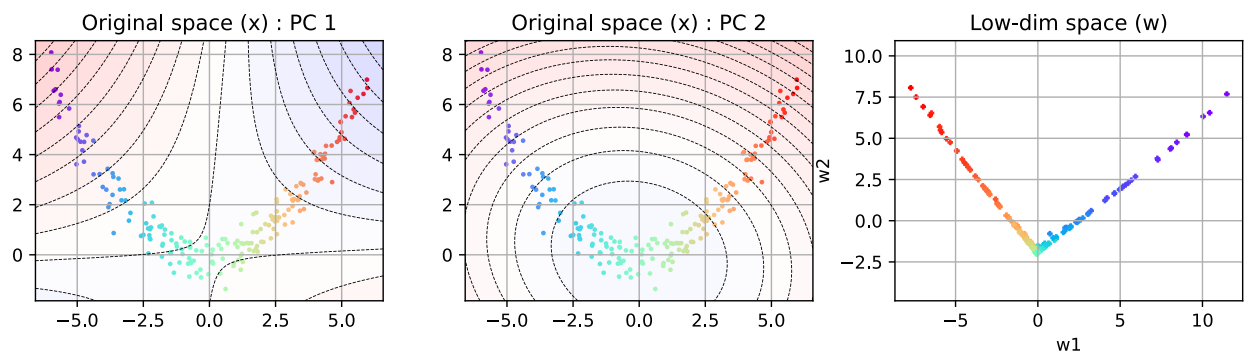
```
plt.figure(figsize=(8,4))
plot_kpca(X, W, kpca, showcontours=True, Y=Y)
```



- Example: 2 PCs
 - 2nd PC corresponds to the distance from the center

```
In [9]: # run KPCA
kpca = decomposition.KernelPCA(n_components=2, kernel='poly', gamma=0.15, degree=2, coef0=0)
W = kpca.fit_transform(X)

plt.figure(figsize=(10,3))
plot_kpca(X, W, kpca, showcontours=True, Y=Y)
```



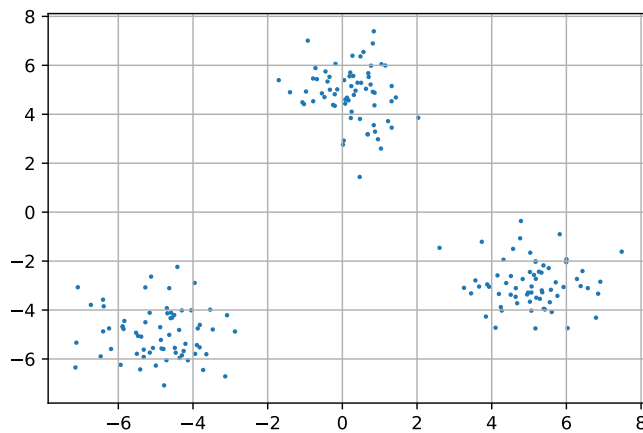
RBF kernel

- principal components separate the data into clusters
- coefficient is distance to clusters

Example

- data with 3 clusters

```
In [11]: plt.scatter(X[:,0], X[:,1], s=6, edgecolor="none")
plt.grid(True);
```

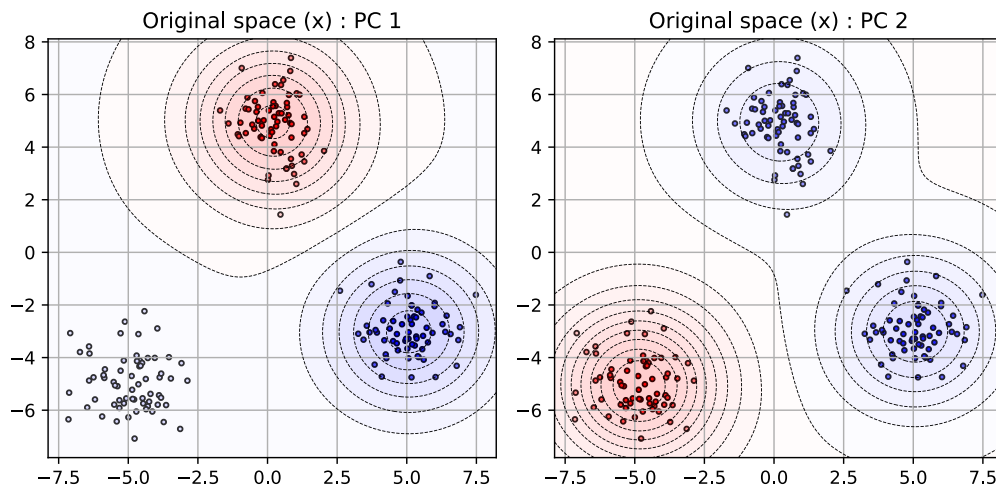


- The first 2 PCs can split the data into 3 clusters
 - the color of the datapoint corresponds to the coefficient value.

In [12]:

```
# run KPCA
kpca = decomposition.KernelPCA(n_components=8, kernel='rbf', gamma=0.15)
W = kpca.fit_transform(X)

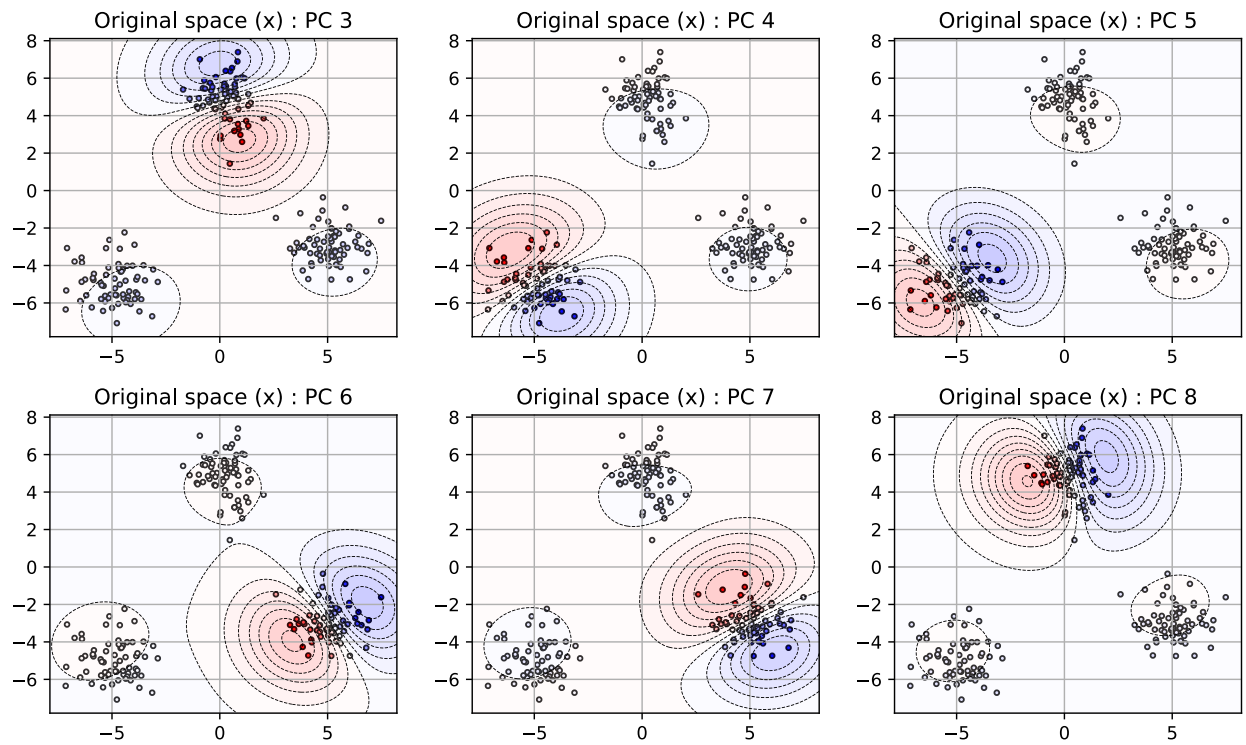
plt.figure(figsize=(8,4))
plot_kpca(X, W, kpca, showcontours=True, showpcs=[0,1], colorcoefs=True)
```



- The remaining 6 PCs split each cluster into halves
 - multiple splits in orthogonal directions

In [13]:

```
plt.figure(figsize=(10,6))
plot_kpca(X, W, kpca, showcontours=True, showpcs=range(2,8), colorcoefs=True)
```



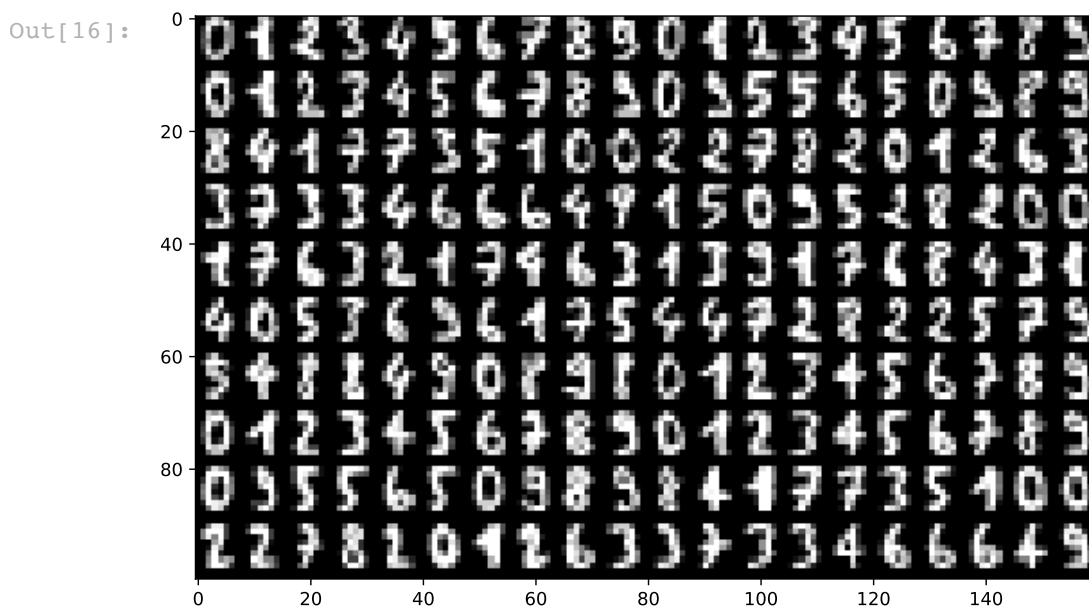
Example on digit images

- 8 x 8 images -> 64D vector

```
In [14]: digits = datasets.load_digits(n_class=10)
X = digits.data
Y = digits.target

# randomly split data into training and testing
trainX, testX, trainY, testY = \
    model_selection.train_test_split(X, Y,
    train_size=0.8, test_size=0.2, random_state=4487)
```

```
In [16]: digitfig
```

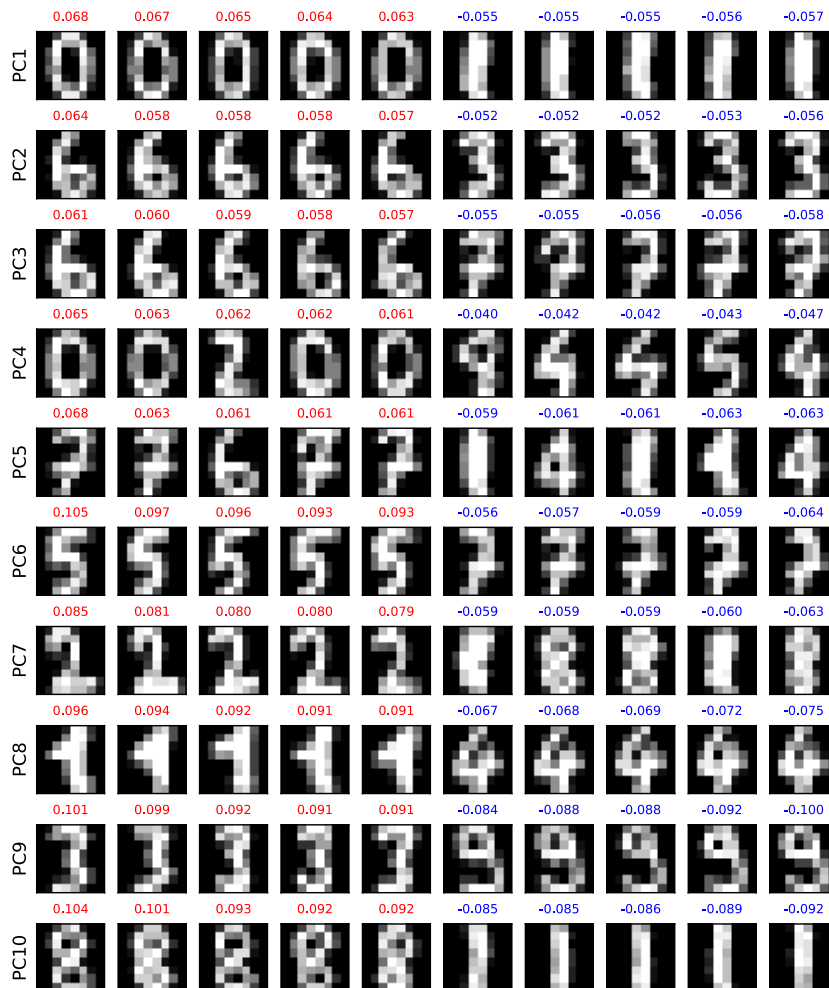


- Apply KPCA with RBF kernel
 - (parallelize with `n_jobs`)

```
In [17]: kpca = decomposition.KernelPCA(n_components=10, kernel='rbf', gamma=0.001, n_jobs=-1)
         trainW = kpca.fit_transform(trainX)
```

- Top-5 positive and negative prototypes for each PC
 - the number is the a_i value for that image.
 - from the prototypes, the PCs are modeling the differences in appearance between digits

```
In [19]: plt.figure(figsize=(8,10))
         plot_kbasis(kpca, (8,8), trainX)
```



Classification experiment

- use KPCA coefficients as the new representation
 - train a logistic regression classifier
 - try different numbers of components
 - Note: can do this efficiently by selecting a subset of KPCA components.

```
In [20]: # apply kernel PCA
         kpca = decomposition.KernelPCA(n_components=60, kernel='rbf', gamma=0.001, n_jobs=-1)
         trainW = kpca.fit_transform(trainX)
         testW = kpca.transform(testX)
```

```
In [21]: ncs = [5, 10, 15, 20, 25, 30, 35, 40, 45, 50, 55, 60]
         accs = []
         for nc in ncs:
             # extract the first nc PCs
             trainWnew = trainW[:,0:nc]
             testWnew = testW[:,0:nc]

             # train classifier
```

```

logreg = linear_model.LogisticRegressionCV(Cs=logspace(-4,4,10), cv=5, n_jobs=-1, max_iter=5000)
logreg.fit(trainWnew, trainY)

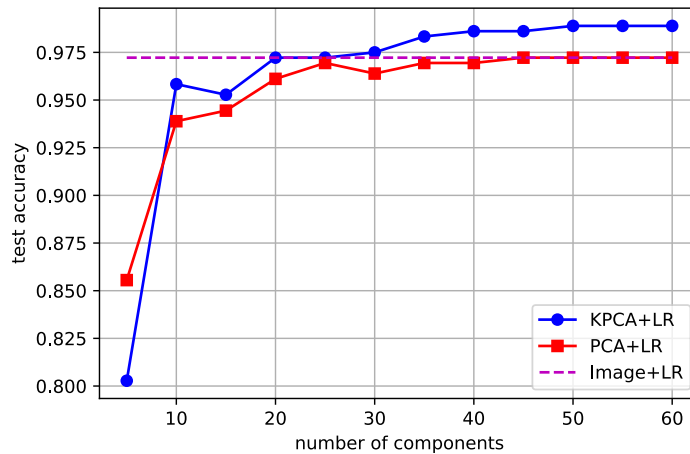
# test classifier
predYtest = logreg.predict(testWnew)
acc = metrics.accuracy_score(testY, predYtest)
accs.append(acc)

```

- Classification results on test set
 - KPCA can improve the performance, compared with PCA and raw image.

In [25]: kfig

Out[25]:



KPCA Summary

- Use kernel trick to perform PCA in high-dimensional space.
 - Coefficients are based on a non-linear projection of the data.
 - The type of projection is based on the kernel function selected.
- Using RBF kernel, KPCA can split the data into clusters.