

Unsupervised Learning

- No targets
- Why use it ?
 - Understand your features
 - Better use of features in supervised models

Plan

- Principal Components
 - Highly popular model for dimensionality reduction
- Clustering
 - K-means to cluster samples
 - Hierarchical clustering
- Recommender systems
 - Netflix prize
 - Pseudo SVD

Alternate basis

We can find an *alternate* set of n basis vectors of length n

$$\tilde{\mathbf{v}}_{(1)}, \dots, \tilde{\mathbf{v}}_{(n)}$$

and translate $\mathbf{x}^{(i)}$ into coordinates $\tilde{\mathbf{x}}^{(i)}$ in the alternate basis

$$\mathbf{x}^{(i)} = \sum_{j=1}^n \tilde{\mathbf{x}}_j^{(i)} * \tilde{\mathbf{v}}_{(j)}$$

PCA finds alternate basis $\tilde{\mathbf{v}}_{(1)}, \dots, \tilde{\mathbf{v}}_{(n)}$

- $\tilde{\mathbf{v}}_{(j)}$ is called *Principal Component j*
- That are mutually orthogonal
$$\tilde{\mathbf{v}}_{(j)} \cdot \tilde{\mathbf{v}}_{(j')} = 0, \text{ for } j \neq j'$$
- $\tilde{\mathbf{v}}_{(j)}$ has more variation than $\tilde{\mathbf{v}}_{(j')}$ for $j' > j$,

The number of basis vectors in the original and alternate bases is both n .

Suppose we reduced the number of alternate basis vectors to $r \leq n$.

- We set $\tilde{\mathbf{x}}_j^{(i)} = 0$ for $j > r$

This is the *reduced dimension* approximation of $\mathbf{x}^{(i)}$.

$$\mathbf{x}^{(i)} = \sum_{j=1}^r \tilde{\mathbf{x}}_j^{(i)} * \tilde{\mathbf{v}}_{(j)}$$

Since the basis vectors are ordered such that $\tilde{\mathbf{v}}_{(j)}$ captures more variation than $\tilde{\mathbf{v}}_{(j')}$ for $j' > j$

- Dropping the alternate bases of higher index loss minimal information

So PCA is the process of

- Finding alternate bases $\tilde{\mathbf{v}}$
- The alternate bases capture correlation among original features \mathbf{x}
- Projecting $\mathbf{x}^{(i)}$ onto $\tilde{\mathbf{v}}$ to obtain transformed vector $\tilde{\mathbf{x}}^{(i)}$ of synthetic features
- Choosing an r so that $\tilde{\mathbf{x}}^{(i)}$ is of dimension $r \leq n$



What is PCA

- a way to reduce dimensionality of features
 - Do we really need all 784 pixels in MNIST ?
 - Hedging in Fixed Income
- a way to cluster samples based on similarity of *features*
- contrast this to Decision Tree
 - also clusters samples with similar features
 - but guided by the clusters having same targets

PCA: High Level

TL;DR

- PCA is a technique for creating "synthetic features" from the original set of features
- The synthetic features may better reveal relationships among original features
- May be able to use reduced set of synthetic features (dimensionality reduction)
- Synthetic features as a means of clustering samples
- **All features need (and will be assumed to be) centered: zero mean**
- PCA is very scale sensitive; often normalize each feature to put on same scale

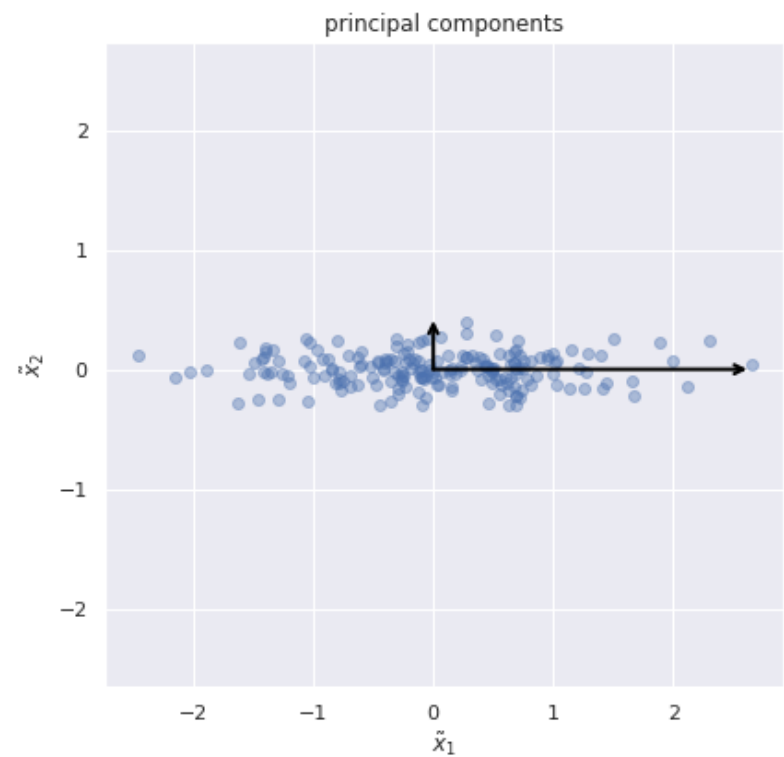
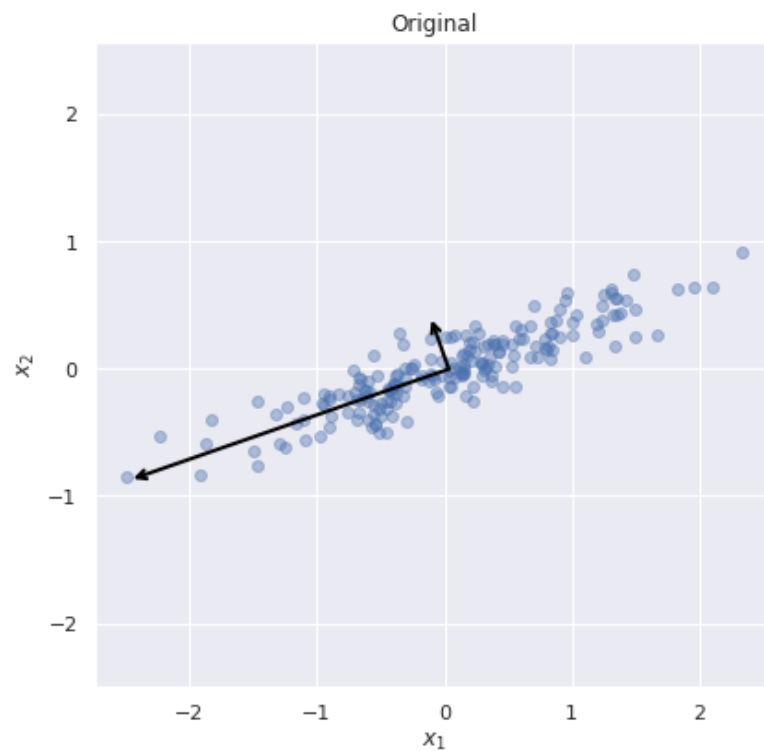
The key ideas behind PCA:

- Dimension reduction: from n "original" features to $n' \leq n$ "synthetic" features
 - synthetic features are more like "concepts" than attributes
 - commonality of purpose
- Feature transformation: from correlated original features to *independent* synthetic features
- Order of "importance" of synthetic features
 - Drop the less important synthetic features: dimension reduction
 - Important relative to feature variation *not* prediction
 - unlike Decision Trees; there is no target

Preview

In one picture:

```
In [4]: X = vp.create_data()  
vp.show_2D(X)
```



The points in the left and right plots are the same, except for the coordinate system.

- Left plot: coordinate system is the horizontal and vertical axes, as usual
- Right plot: coordinate system is the dark, arrowed lines -- identical to the lines on the left plot

Bottom line: same data, expressed in a different way

- Left plot: features $\mathbf{x}_1, \mathbf{x}_2$
- Right plot: features $\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2$

In the left plot, we can clearly see that the data set's features $\mathbf{x}_1, \mathbf{x}_2$ are correlated.

In the right plot: $\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2$ are

- independent
- with $\tilde{\mathbf{x}}_1$ expressed greater variation

The arrowed vectors in the original

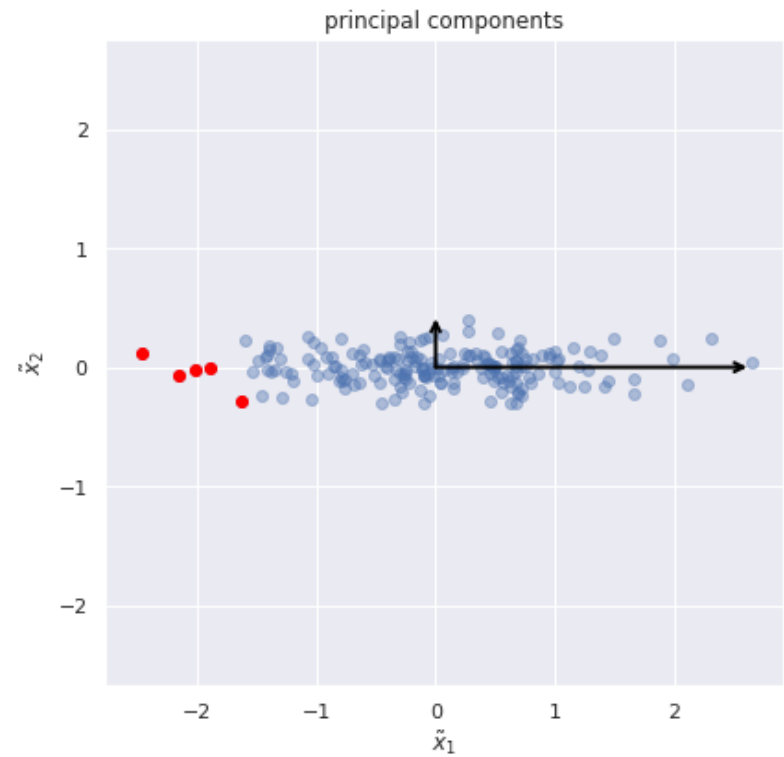
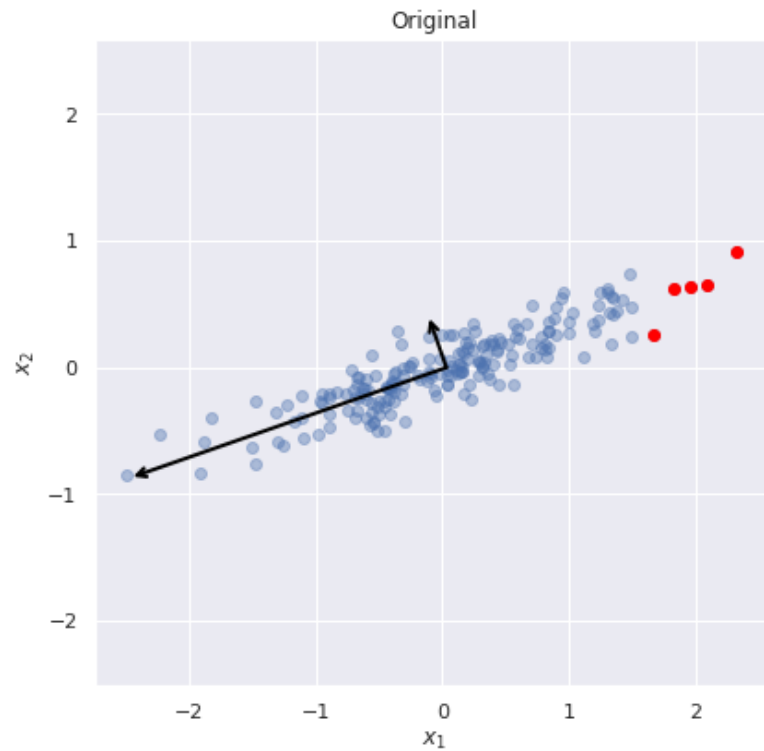
- correspond to the components: the alternate basis vectors
- the first component is in the negative direction
- so the plot on the right, which has this component in the *positive* direction, is "flipped"
 - we illustrate this by showing the mapping of the points in red below

You can see that the alternate basis

- First basis vector (component) points left
- So the first coordinate (feature) in the alternate basis is opposite the first coordinate in the original

We can more easily see it by observing how the red examples behave

```
In [5]: vp.show_2D(X, points=X[ X[:,0] > 1.5 ])
```



To summarize

$$\mathbf{X} = \tilde{\mathbf{X}}V^T$$

That is

- Examples \mathbf{X}
 - original features $\mathbf{x}^{(i)}$
- Re-expressed in new basis V^T
 - synthetic features $\tilde{\mathbf{x}}^{(i)}$
 - new basis $\setminus \mathbf{V}^T$ may have fewer dimensions $r \leq n$ than original



PCA via Matrix factorization

$$\mathbf{X} = \{\mathbf{x}^{(i)} \mid i = 1, \dots, m\}$$

No targets !

\mathbf{X} is **zero centered** (subtract mean from each feature)

Our goal is to

- Find $\tilde{\mathbf{X}}, V^T$
- Such that
$$\mathbf{X} = \tilde{\mathbf{X}}V^T$$

Decomposing \mathbf{X} into a product (as above) is called *matrix factorization*

Some types of matrix factorization we'll mention

- Singular Value Decomposition
- Eigen Decomposition
- CUR Decomposition

We will first derive PCA using Singular Value Decomposition of \mathbf{X}

Singular Value Decomposition (SVD) Factorization

Factor matrix X into product of 3 matrices:

$$\mathbf{X} = U\Sigma V^T$$

- U : $m \times n$, columns are orthogonal unit vectors
 - $UU^T = I$
- Σ : $n \times n$ diagonal matrix
 - $\text{diag}(\Sigma) = [\sigma_1, \sigma_2, \dots, \sigma_n]$
- V : $n \times n$, columns are orthogonal unit vectors
 - $VV^T = I$

Moreover, the diagonal elements of Σ are in descending order of magnitude

$$\sigma_j > \sigma_{j'}$$

for $j' > j$

Let V^T be the new basis vectors

- Since $VV^T = I$ these vectors are orthogonal

We need to find the synthetic features \mathbf{X} relative to these bases

$$\mathbf{X} = \tilde{\mathbf{X}}V^T$$

$$\mathbf{X} = \tilde{\mathbf{X}}V^T$$

$$U\Sigma V^T = \tilde{\mathbf{X}}V^T \quad \text{factorization } \mathbf{X} = U\Sigma V^T$$

$$U\Sigma V^T V = \tilde{\mathbf{X}}V^T V \quad \text{multiple both sides by } V$$

$$U\Sigma = \tilde{\mathbf{X}} \quad \text{since } V^T V = I$$

Thus SVD gives us both $\tilde{\mathbf{X}}$ and V^T as desired.

In fact since

$$\tilde{\mathbf{X}} = U\Sigma$$

there is interesting structure in $\tilde{\mathbf{X}}$

Recall that U is orthonormal

$$UU^T = I$$

so that its vectors are unit length

Since Σ is diagonal (all zero for non-diagonal elements)

- The diagonal of Σ , denoted

$$\text{diag}(\Sigma) = [\sigma_1, \sigma_2, \dots, \sigma_n]$$

scales the columns of U

$$\begin{aligned}(U\Sigma)^{(i)} &= U^{(i)} * \text{diag}(\Sigma) \\ &= [U_1^{(i)} * \sigma_1, U_2^{(i)} * \sigma_2, \dots, U_n^{(i)} * \sigma_n]\end{aligned}$$

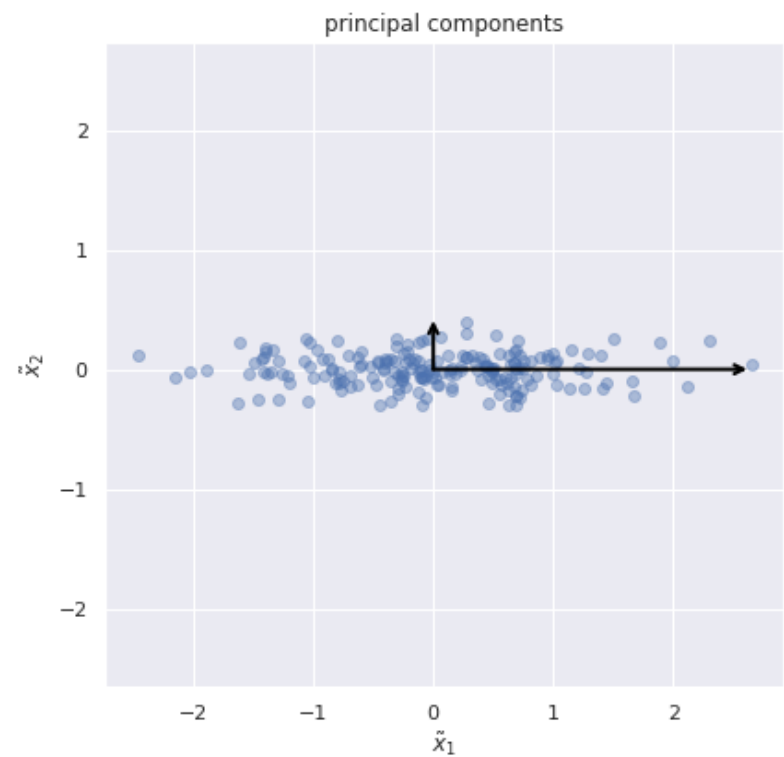
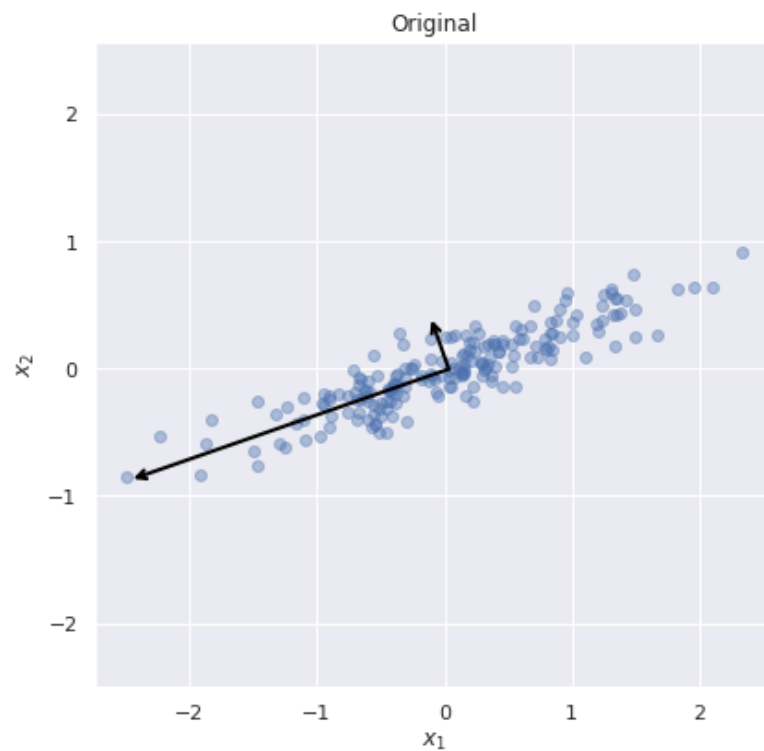
Thus

- U is a "standardized" version of features $\tilde{\mathbf{X}}$
 - unit standard deviation
- $\tilde{\mathbf{X}} = U\Sigma$ is the non-standardized features

A picture may clarify the distinction between the standardized and non-standardized $\tilde{\mathbf{X}}$.

Here is the non-standardized $\tilde{\mathbf{X}}$ that we've seen previously

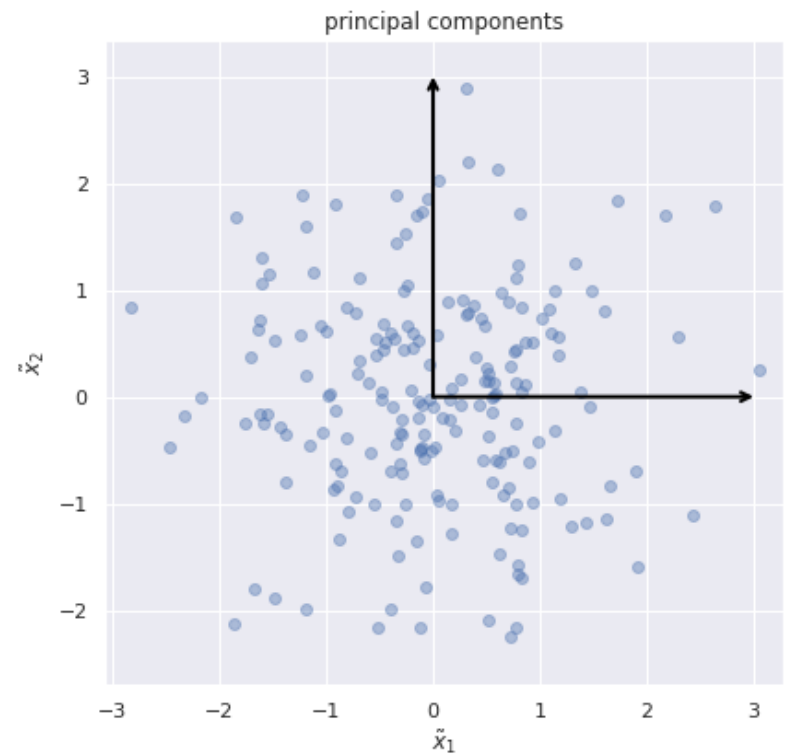
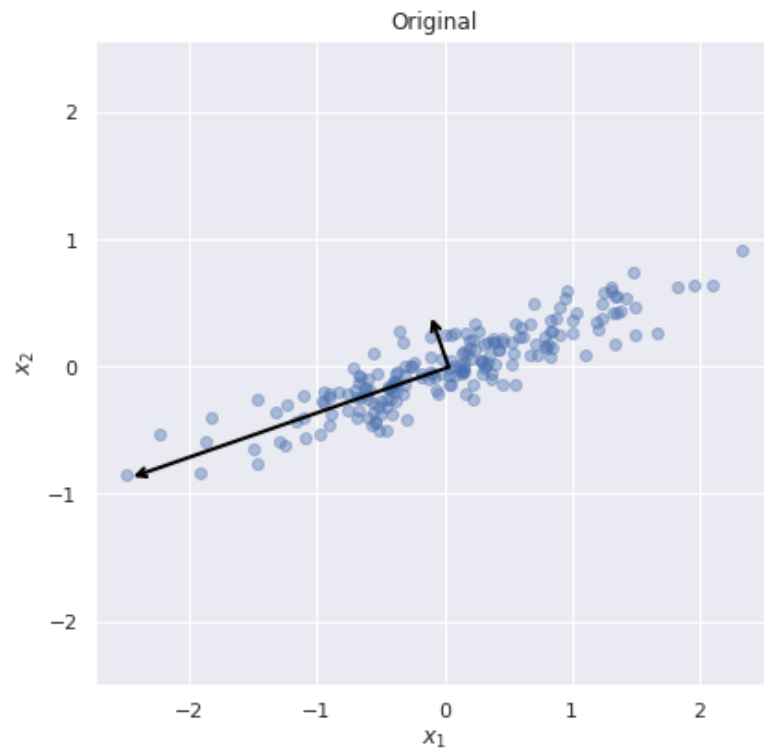
```
In [6]: X = vp.create_data()  
vp.show_2D(X)
```



And here is the standardized plot

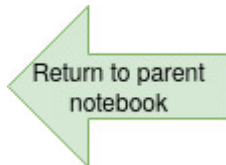
- The length of each basis vector is 1
- Rather than σ_j
- By stretching each component by σ_j we recover the non-standardized plot

```
In [7]: vp.show_2D(X, whiten=True)
```



So

- The synthetic features $\tilde{\mathbf{X}}$
- Are "standardized" synthetic features U
- Scaled by $\text{diag}(\Sigma)$



Dimensionality reduction

Thus far we have

- exactly replicated \mathbf{X} via new bases V^T
- $$\mathbf{X} = \tilde{\mathbf{X}}V^T$$

But the length of examples in \mathbf{X} and $\tilde{\mathbf{X}}$ are both n , so the number of features is unchanged.

We will now reduce the dimension of $\tilde{\mathbf{X}}$ to $r \leq n$

Suppose we set

$$\sigma_j = 0$$

for all $j > r$

Then $\tilde{\mathbf{X}}_j^{(i)} = 0$ for $j > r$ since the higher order features have been scaled by 0.

Effectively we have reduced the dimensionality of $\tilde{\mathbf{X}}$ from n to $r \leq n$.

By eliminating columns that are all 0 in U , Σ and V^T we get an *approximation* \mathbf{X}' of \mathbf{X}

$$\mathbf{X}' \approx \mathbf{X}$$

$$\mathbf{X}' = U' \Sigma' (V^T)'$$

where

Note that

- \mathbf{X}' is reduced dimension $r \leq n$
- $(V^T)'$ consists of
 - r rows each of length n
 - each row is a new basis vector

Best lower rank approximation of \mathbf{X}

\mathbf{X}' is an approximation of \mathbf{X} .

We define what the approximation error

$$||\mathbf{X}' - \mathbf{X}||_2 = \sum_{i,j} (\mathbf{X}'_{i,j} - \mathbf{X}_{i,j})^2$$

The above is called the Froebenius Norm (and looks like MSE in form).

The "best" approximation \mathbf{X}' of \mathbf{X} is the one with the smallest norm.

Let

$$\mathbf{X} = U\Sigma V^T$$

be the SVD of \mathbf{X} where

- Σ is dimension $(n \times n)$ with diagonal element $\sigma_1, \dots, \sigma_n$ such that
- $\sigma_i > \sigma_j$ for $i < j$

$\mathbf{X}_{\text{d-proj}}$ is the best rank d approximation of \mathbf{X} , i.e,

$$\mathbf{X}_{\text{d-proj}} = \operatorname{argmin}_B \|\mathbf{X} - B\|_2$$

when

$$\mathbf{X}_{\text{d-proj}} = USV^T$$

where

S is the $(n \times n)$ diagonal matrix with

- first d diagonal elements $\sigma_1, \dots, \sigma_d$
- remaining $(n - d)$ diagonal elements set to 0

S is a "truncated" version of the Σ from SVD (setting diagonal members to 0 for $j > d$)

The zero elements of S render the last $(n - d)$ columns of U and rows of V^T irrelevant.

$$\mathbf{X}_{\text{d-proj}} = U_{\text{d-proj}} \Sigma_{\text{d-proj}} V_{\text{d-proj}}^T$$

In words, this effectively reduces the dimensions of U, Σ, V^T :

- $U_{\text{d-proj}}$ to dimension (m, d)
- $\Sigma_{\text{d-proj}}$ to (d, d)
- $V_{\text{d-proj}}^T$ to (d, n)

In effect, we retain only the first d synthetic features.

So we have reduced the dimension of \mathbf{X} from n features to $\tilde{\mathbf{X}}$ with d synthetic features.

Why dropping low significance synthetic features is low error

Let's examine the reconstruction of \mathbf{X} from $U\Sigma V^T$

$$\mathbf{X} = U\Sigma V^T$$

$$\mathbf{X}_j^{(i)} = U^{(i)} \cdot (\Sigma V^T)_j \quad \text{definition of matrix multiplication}$$

$$= U^{(i)} \cdot (\Sigma T)_j \quad \text{Let } T = V^T$$

$$= U^{(i)} \cdot \begin{pmatrix} T_j^{(1)} * \sigma_1 \\ T_j^{(2)} * \sigma_2 \\ \vdots \\ T_j^{(n)} * \sigma_n \end{pmatrix} \quad \Sigma \text{ scales the rows of } V^T$$

$$= \sum_{k=1}^n U_k^{(i)} * T_j^{(k)} * \sigma_k \quad \text{definition of dot product}$$

That is:

- Example i , feature j
 - is the product of n terms
 - where the k^{th} term is scaled by σ_k

- Keeping *only* the first d terms of $\mathbf{X}_j^{(i)}$

- results in error of:

$$\sum_{k>d}^n U_k^{(i)} * T_j^{(k)} * \sigma_k$$

- recall that $\sigma_j > \sigma_{j'}$ for $j < j'$
- so the $(n - d + 1)$ dropped terms are scaled by the *smallest* σ_k
 - if $\sigma_k \approx 0$, for $k > d$
 - the approximation error is minimized

That's why we drop the higher order synthetic features: their contribution to $\mathbf{X}_j^{(i)}$ is smallest

How many dimensions to keep ?

Since the diagonal elements of Σ are ordered

- We keep a cumulative sum of σ^2 , which will sum to 1.
- Choose to keep the first k synthetic features
 - where the cumulative sum up to (and including) k is greater than some fraction
 - e.g., 95%.

We will illustrate this in the following example.



The inverse transformation

We have shown how to transform original features \mathbf{X} to synthetic features $\tilde{\mathbf{X}}$.

How about inverting the transformation: recover \mathbf{X} from $\tilde{\mathbf{X}}$?

Since

$$\mathbf{X} = \tilde{\mathbf{X}}V^T \quad \text{definition}$$

$$\mathbf{X}V = \tilde{\mathbf{X}}V^TV \quad \text{multiply both sides by } V$$

$$\mathbf{X}V = \tilde{\mathbf{X}} \quad \text{since } V^TV = I$$

So

- V transforms from original features \mathbf{X} to synthetic feature $\tilde{\mathbf{X}}$
- V^T transforms synthetic features $\tilde{\mathbf{X}}$ to original features \mathbf{X}



Other decompositions of \mathbf{X}

Eigen decomposition of covariance matrix of \mathbf{X}

There is another matrix factorization method known as Eigen Decomposition.

Eigen decomposition, unlike SVD, only works on symmetric matrices M :

$$M = W\Lambda W^T$$

where $WW^T = I$

We can obtain the PCA from the Eigen Decomposition of $\mathbf{X}\mathbf{X}^T$

- the covariance matrix of \mathbf{X} (i.e., original feature covariance)
- the covariance matrix is symmetric, as required

We can relate the SVD of \mathbf{X} to the Eigen decomposition of $\mathbf{X}\mathbf{X}^T$ as follows:

$$\begin{aligned}\mathbf{X}^T \mathbf{X} &= \mathbf{V} \Sigma \mathbf{U}^T \mathbf{U} \Sigma \mathbf{V}^T && \text{from SVD } \mathbf{X} = \mathbf{U} \Sigma \mathbf{V}^T \\ \mathbf{X}^T \mathbf{X} &= \mathbf{V} \Sigma \Sigma^T \mathbf{V}^T && \text{since } \mathbf{U}^T \mathbf{U} = \mathbf{I}\end{aligned}$$

Similarly, we can show

$$\mathbf{X}\mathbf{X}^T = \mathbf{U} \Sigma \Sigma^T \mathbf{U}^T \text{ since } \mathbf{V} \mathbf{V}^T = \mathbf{I}$$

Setting

- $\Lambda = \Sigma \Sigma^T$
- $\mathbf{W} = \mathbf{U} = \mathbf{V}$ we get $\mathbf{X} = \mathbf{W} \Lambda \mathbf{W}^T$, the Eigen Decomposition of $\mathbf{X}\mathbf{X}^T$.

The V that transforms \mathbf{X} (original features) to $\tilde{\mathbf{X}} = \mathbf{X}V$ (synthetic features)

- Can be computed directly from SVD
- Or by creating covariance matrix $\mathbf{X}\mathbf{X}^T$ and using Eigen decomposition.

SVD is more commonly used

- There are many fast implementations of SVD
- There is no need to compute the big covariance matrix $\mathbf{X}\mathbf{X}^T$

Other factorization methods

- CUR method

$$\text{CUR}(A) = C \cdot U \cdot R$$

- C chosen from Columns of R
- R chosen from Rows of A

Example: Reconstructing \mathbf{x} from $\tilde{\mathbf{x}}$ and the principal components

It may be helpful to visualize

- The transformation from example features $\mathbf{x}^{(i)}$
- To synthetic features $\tilde{\mathbf{x}}^{(i)}$

We will use a subset of the "smaller digits" (8×8) data

```
In [8]: subset1 = [ 0, 4, 7, 9 ]  
rh_digits = unsupervised_helper.Reconstruct_Helper( subset=[])  
rh_digits.create_data_digits(subset=subset1)
```

```
In [9]: n_components = 8  
        _ = rh_digits.fit(n_components=n_components)
```

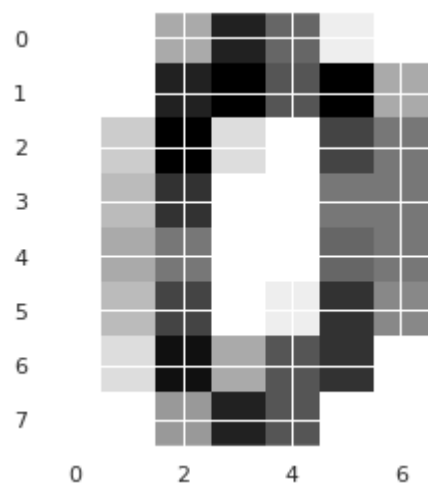


```
In [10]: # Which example to show  
data_idx = 0  
fig0, ax0, figm, axm, figc, axc = rh_digits.show_data_comp(data_idx=data_idx)  
  
plt.close(fig0)  
plt.close(figm)  
plt.close(figc)
```

Here is one example:

In [11]: fig0

Out[11]:

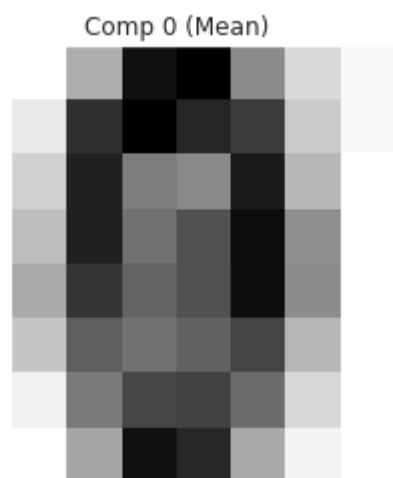


And here is the "mean" of \mathbf{X}

- Recall: we assume \mathbf{X} has been zero-centered
- So the mean is subtracted before performing PCA
- Which means it has to be *added* to the reconstructed image

In [12]: figm

Out[12]:

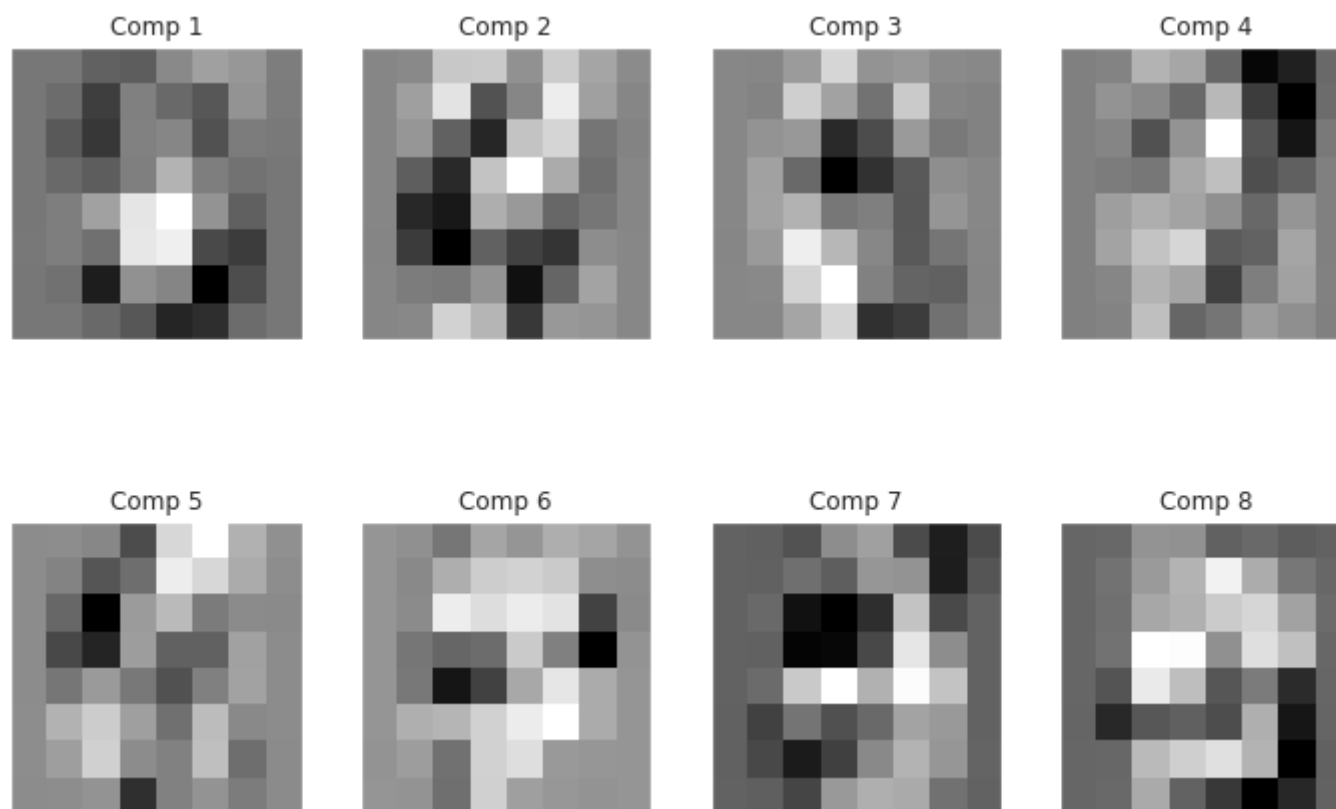


And here are the Principal Components (new bases)

- We performed PCA with a reduced number of components ($r < n$)
- There are `n_components` such basis vector
- *Each* component is of length n , but there are only $r < n$ components

In [13]: figc

Out[13]:



It's not necessarily easy to interpret the components, particularly when n is large

- Component 1 *might* be the "concept" corresponding to the digit 0
- Component 4 *might* be the "concept" corresponding to the digit 4
- The other components *might* be *partial* shape concept, rather than entire digits

Let's progressively examine

$$\mathbf{x}^{(i)} = \sum_{j=1}^r \tilde{\mathbf{x}}_j^{(i)} * (V^T)^{(j)}$$

It may be helpful to remind ourselves of the shapes of each element in the equation

- $\mathbf{x}^{(i)}$ is of length n
- The components V^T are $(r \times n)$
 - Each component is of length n
 - There are $r < n$ components (e.g., $r = \text{n_components}$)
- $\tilde{\mathbf{x}}^{(i)}$ is of length r

Note

- We treat the "mean\X\$ as component 0
- With weight 1

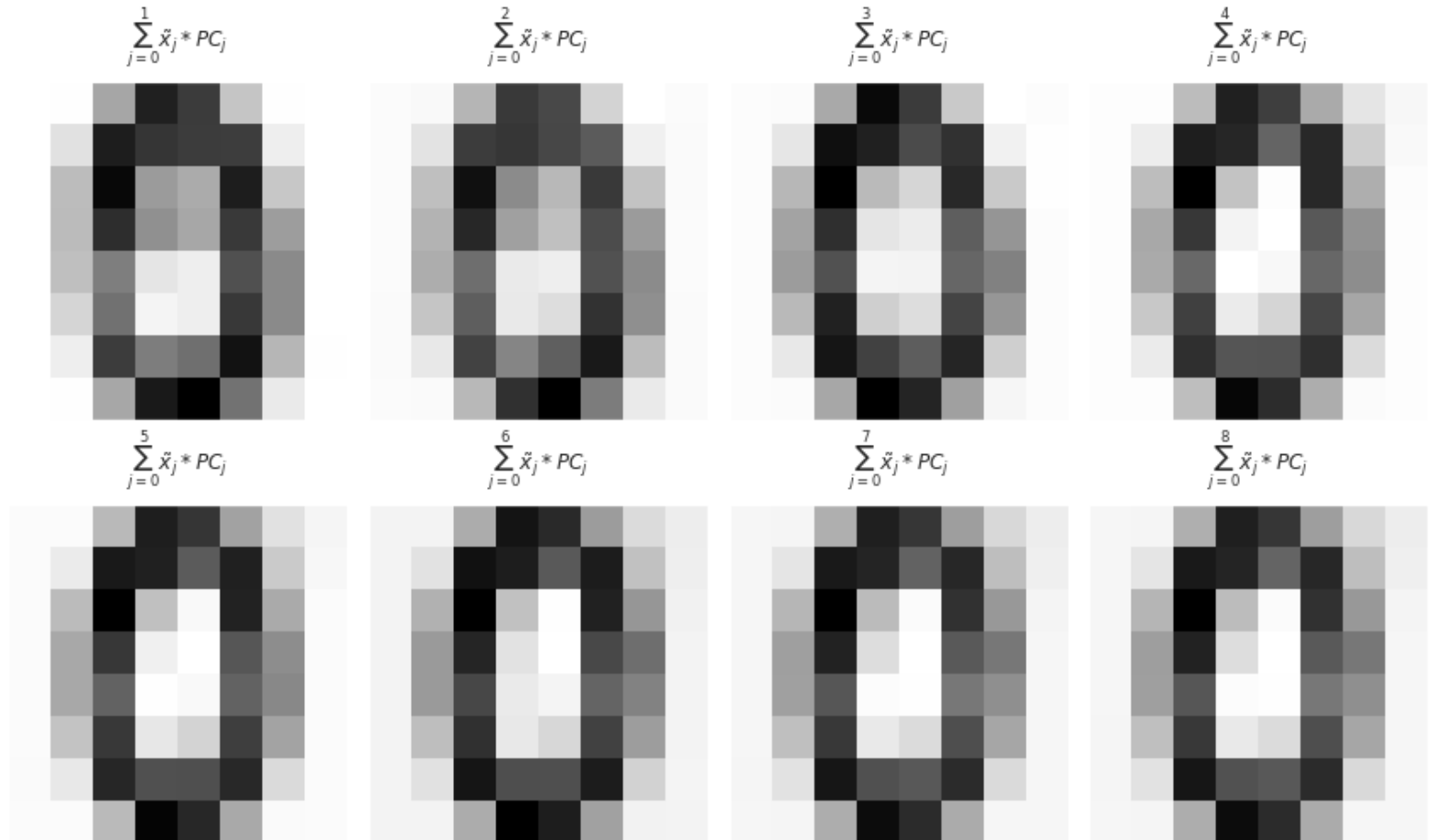
So we construct an approximation of $\mathbf{x}^{(i)}$

- By adding weighted components $(V^T)^{(j)}$, each of length n
- The weight associated with component j is $\tilde{\mathbf{x}}_j^{(i)}$
- So the sum is of length n

```
In [14]: figi, axsi, fig_comp, axs_comp, x_tilde, error = rh_digits.show_recon(data_idx=d
ata_idx)
plt.close(figi)
plt.close(fig_comp)
```

```
In [15]: figi
```

Out[15]:



You can see that the approximation using just the first component (and the mean)

- Is already a good approximation of $\mathbf{x}^{(i)}$
- Somewhat confirming our *guess* that component 1 represents the concept 0

We can confirm this by looking at $\tilde{\mathbf{x}}^{(i)}$

In [16]: `print("x tilde = ", x_tilde)`

```
arg_max = np.argmax(x_tilde)
```

```
print("Largest feature at index {idx:d}".format(idx=arg_max+1))
```

```
x tilde = [ 18.94187266  5.09553028 -11.1173937    6.39761384 -0.96183808
           3.24941375  2.53344604  0.07949071]
```

```
Largest feature at index 1
```

As you can see, the magnitude of $\tilde{\mathbf{x}}_1^{(i)}$ is the largest among $\{\tilde{\mathbf{x}}_j^{(i)} \mid 1 \leq j \leq r\}$

In fact, we might try to confirm our intuition

- By examining $\tilde{\mathbf{x}}^{(i')}$ for all i' where $\mathbf{y}^{(i')} = 0$ (assuming we have targets/labels)

```
In [17]: # Get X tilde and the targets  
Xtilde = rh_digits.dataProj  
y = rh_digits.targets  
  
# Filter to identify examples where target is equal to digit  
digit = 0  
mask = (y == digit)  
Xtilde_digit = Xtilde[mask]
```

```
In [18]: print("x tilde, when y=0:")  
  
for i in range(0,10):  
    print( [ "{x:3.2f}".format(x=x_tilde_j) for x_tilde_j in Xtilde_digit[i] ])
```

x tilde, when y=0:

```
['18.94', '5.10', '-11.12', '6.40', '-0.96', '3.25', '2.53', '0.08']  
['10.94', '11.59', '-8.82', '8.34', '-6.51', '4.21', '4.22', '-4.47']  
['15.16', '8.46', '-9.70', '2.85', '-3.39', '-5.13', '-3.83', '-7.18']  
['21.68', '9.93', '-12.65', '3.27', '1.32', '4.84', '-0.85', '-4.05']  
['13.36', '8.83', '-11.29', '4.24', '-0.35', '-1.16', '0.97', '-12.10']  
['19.57', '7.25', '-9.87', '-3.35', '-1.75', '7.13', '4.02', '0.29']  
['20.21', '9.81', '-7.81', '-2.10', '-2.19', '-1.39', '-0.84', '3.74']  
['10.75', '12.33', '-9.70', '-1.01', '-4.82', '-7.07', '-4.93', '-12.89']  
['17.96', '12.42', '-5.09', '-5.08', '2.47', '-8.21', '-2.90', '-11.12']  
['22.48', '1.51', '-7.25', '-10.33', '3.19', '3.59', '0.19', '4.90']
```

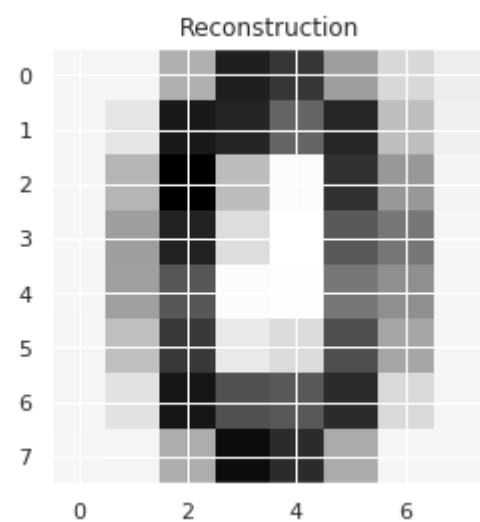
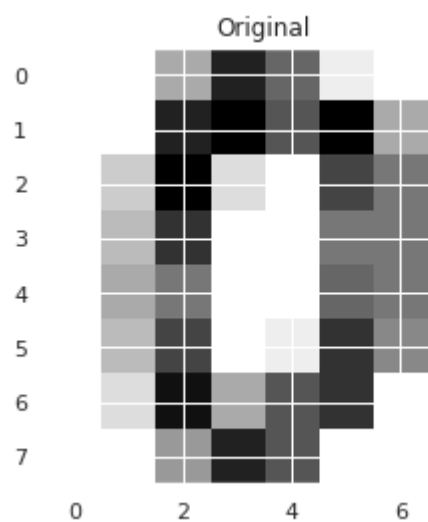
As you can see

- For examples $\mathbf{x}^{(i')}$ where $\mathbf{y}^{(i')} = 0$
- $\tilde{\mathbf{x}}^{(i')}$ is the largest value in $\tilde{\mathbf{x}}^{(i')}$

Here is a comparison of the original $\mathbf{x}^{(i)}$ and its reconstructed approximation

In [19]: fig_comp

Out[19]:



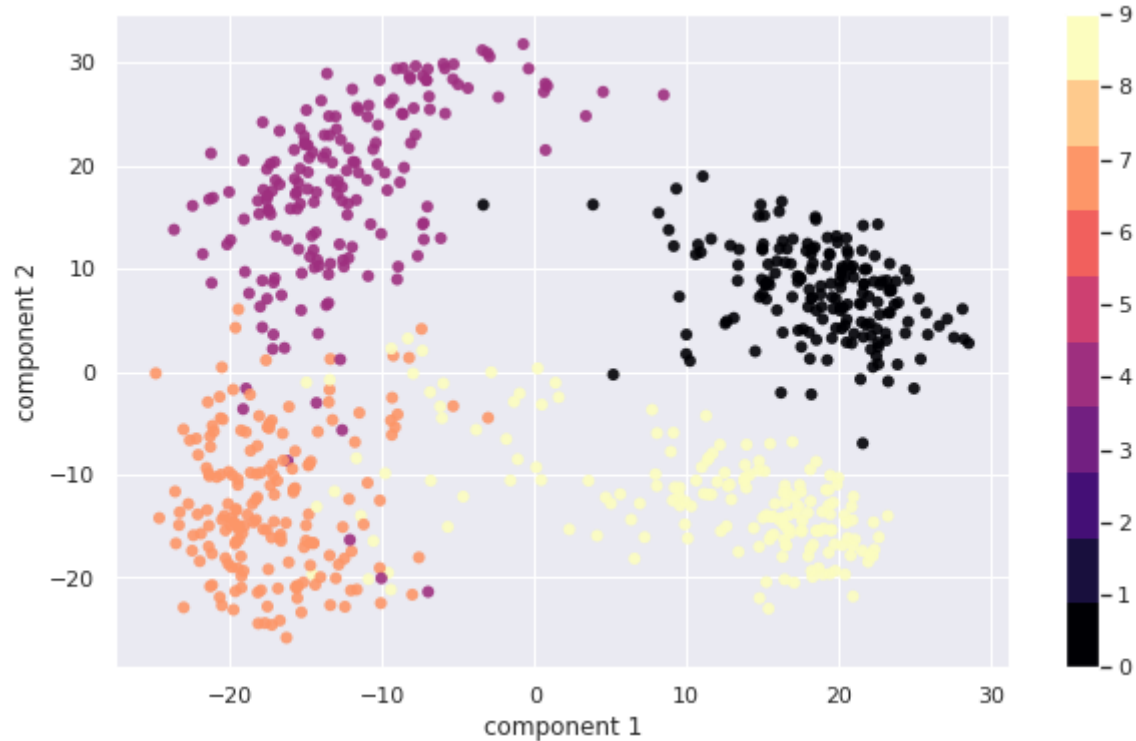
We could try to plot all our examples in r dimensional space

- To see whether examples formed clusters (examples with similar feature vectors of length r)

But $r = \text{n_components}$ is too large in our case; let's just plot using the first two features of $\tilde{\mathbf{x}}$

```
In [20]: vpt = unsupervised_helper.VanderPlas()  
  
print("Number of examples: {n:d}".format(n=Xtilde.shape[0]))  
vpt.digits_subset_show_clustering(Xtilde, y, save_file="/tmp/digits_subset_cluster.jpg" )
```

Number of examples: 718



Since we have targets/labels available (not generally the case for unsupervised learning)

- We can color the points according to their target
- We see that the 4 digits in the restricted examples cluster according to their features in $\tilde{\mathbf{x}}$
- Digit "0" is associated with (high $\tilde{\mathbf{x}}_1$, high $\tilde{\mathbf{x}}_2$)
- Digit "4" is associated with (low $\tilde{\mathbf{x}}_1$, high $\tilde{\mathbf{x}}_2$)



Dimensionality reduction:examples

MNIST example

- 784 features
 - are some redundant ? Can we capture "essence" with fewer pixels ?
 - Consider blocks of black pixels in 4 corners
 - pixel (i,j) highly correlated (across samples) with pixel $(i+1,j)$, $(i-1,j)$, $(i,j+1)$, etc.
 - i.e. in many samples: when pixel (i,j) is black, so are surrounding pixels
 - If we replaced the block with 1 synthetic feature ("block of black in upper left ...")
 - we can reduce number of features (many pixels reduced to single)
 - reconstruction from compressed feature space to original preserves most info

So above goal was to reduce number of dimensions without losing info

Ideally, the reduction would be to a small enough (2 or 3) number of synthetic features

- that we could plot the samples in the transformed synthetic feature space.

Retrieve the full MNIST dataset (70K samples)

We had previously used only a fraction in order to make our demo faster.

```
In [21]: ush = unsupervised_helper.PCA_Helper()  
X_mnist, y_mnist = ush.mnist_init()
```

Retrieving MNIST_784 from cache

```
In [22]: from sklearn.model_selection import train_test_split
X_mnist.shape, y_mnist.shape
X_mnist_train, X_mnist_test, y_mnist_train, y_mnist_test = train_test_split(X_mnist, y_mnist)
X_mnist_train.shape
```

```
Out[22]: ((70000, 784), (70000,))
```

```
Out[22]: (52500, 784)
```

Perform PCA.


```
In [23]: pca_mnist = ush.mnist_PCA(X_mnist_train)
```

```
In [24]: pca_mnist.n_components_  
X_mnist_train_reduced = ush.transform(X_mnist_train, pca_mnist)  
X_mnist_train_reduced.shape
```

```
Out[24]: 154
```

```
Out[24]: (52500, 154)
```

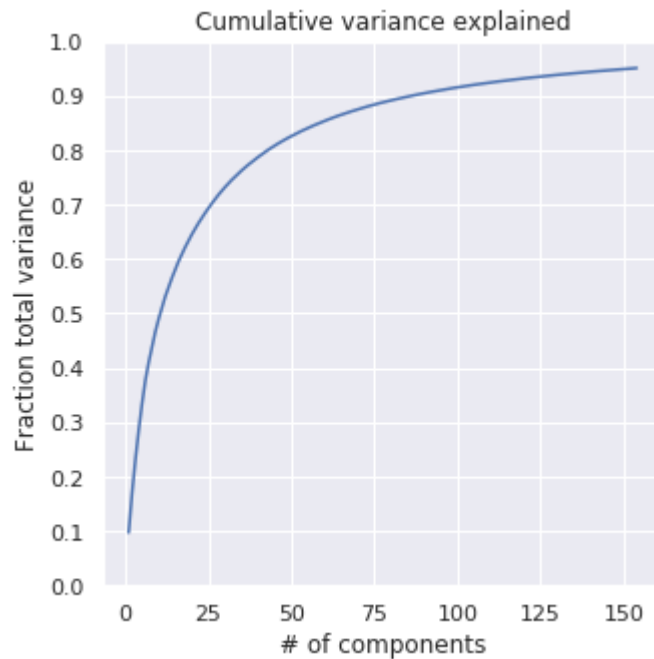
Let's plot the cumulative variance as a function of number of synthetic features.

This can help us determine how many synthetic features to keep.

```
In [25]: _ = ush.plot_cum_variance(pca_mnist)

variance_goal_pct = 95
features_for_goal = ush.num_components_for_cum_variance(pca_mnist, .01 * varianc
e_goal_pct)
print("To capture {f:d}% of variance we need {d:d} synthetic features.".format(f
=variance_goal_pct, d=features_for_goal))
```

To capture 95% of variance we need 154 synthetic features.



So we need only about 20% of the original 784 features to capture 95% of the variance.

We can invert the PCA transformation to go from synthetic feature space back to original features.

That is, we can see what the digits look like when reconstructed from only 154 synthetic features.

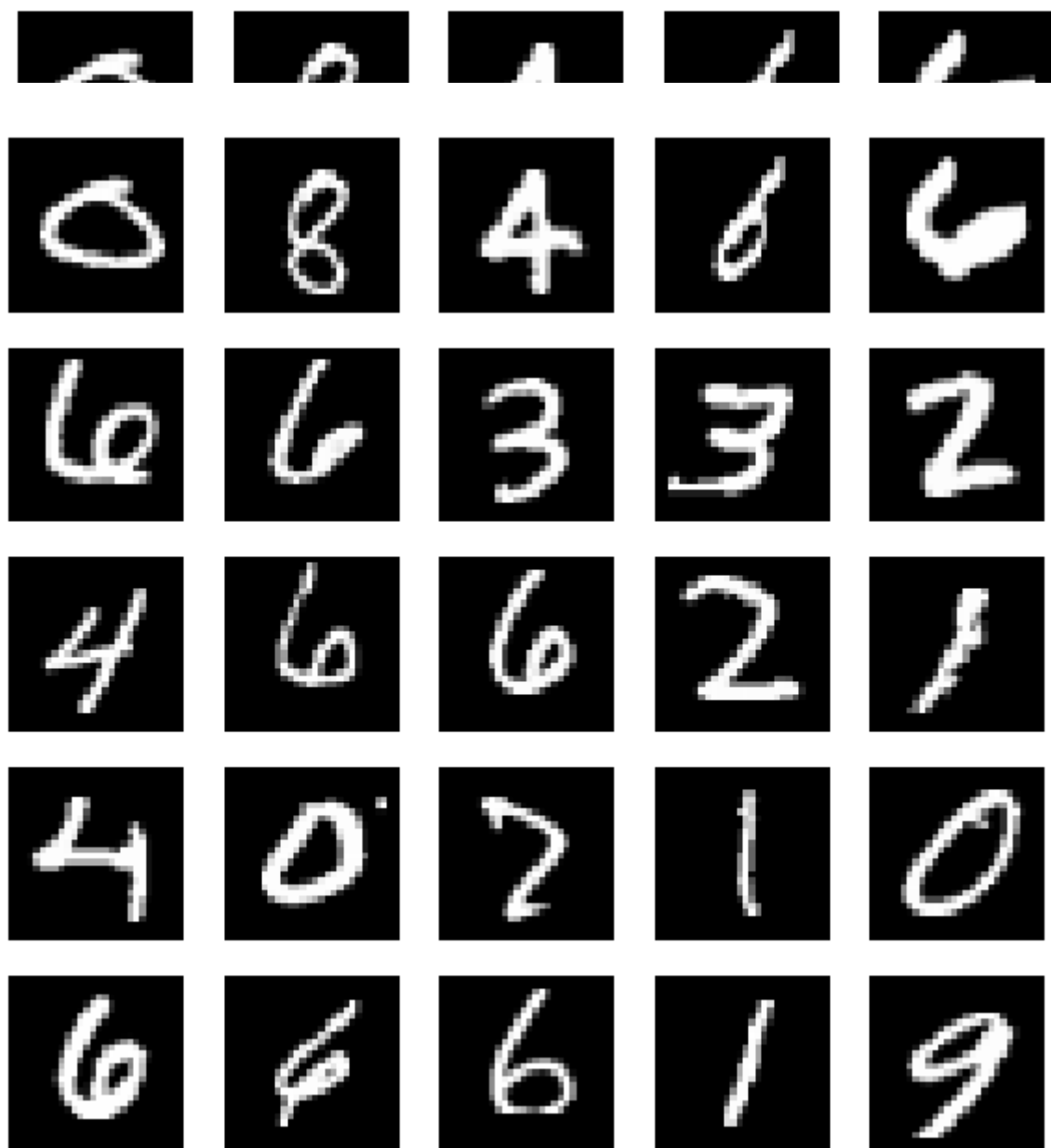
First, let's look at the original:

```
In [26]: X_mnist_train_reduced = ush.transform(X_mnist_train, pca_mnist)
X_mnist_train_reduced.shape

# Show the original dataset
ush.mnh.visualize(X_mnist_train, y_mnist_train)
```

```
Out[26]: (52500, 154)
```

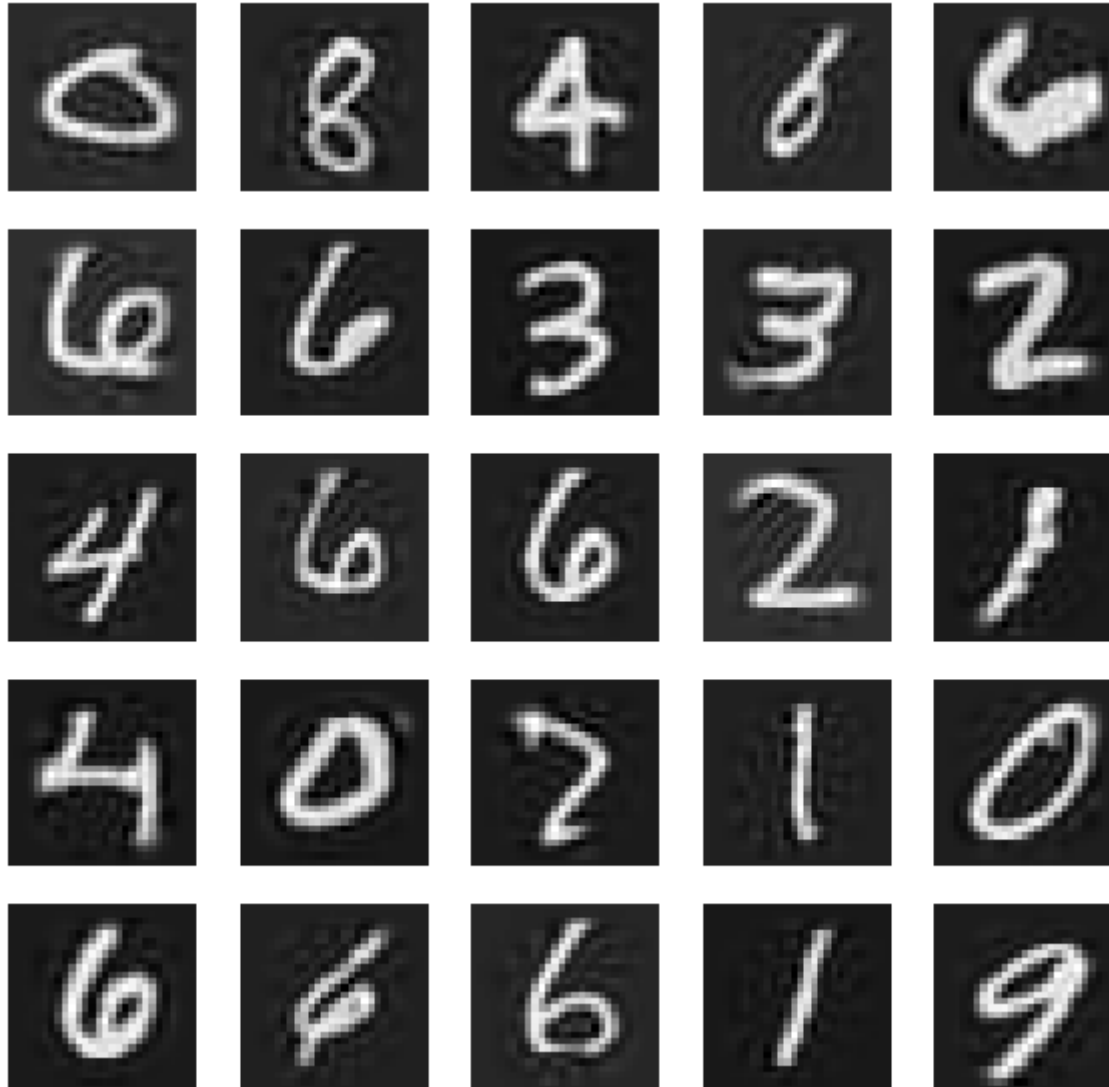
```
Out[26]:
```

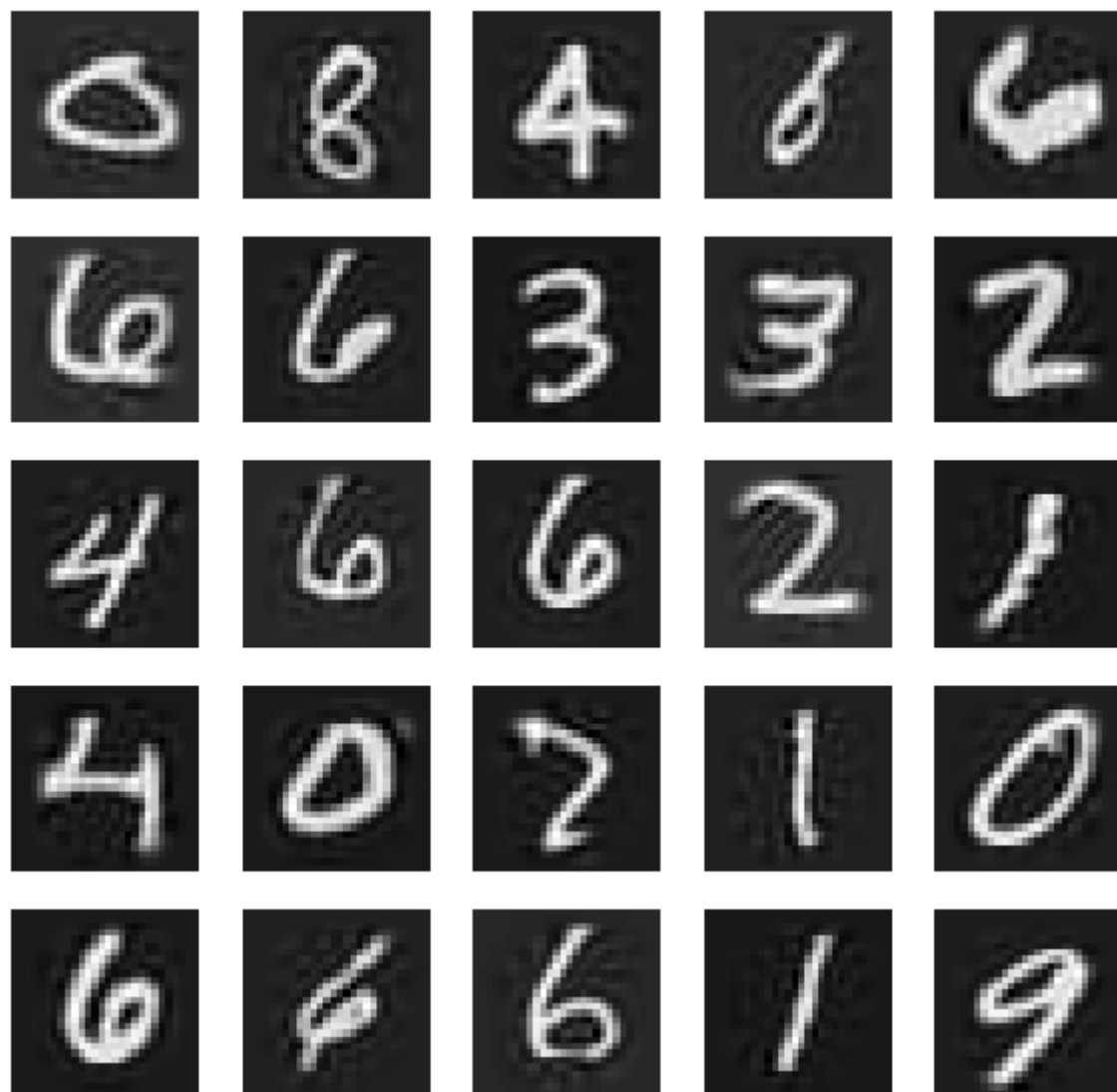


Next, the reconstructed


```
In [27]: X_mnist_train_reconstruct = ush.inverse_transform(X_mnist_train_reduced, pca_mnist)
ush.mnh.visualize(X_mnist_train_reconstruct, y_mnist_train)
```

Out[27]:





A little fuzzy, but pretty good.

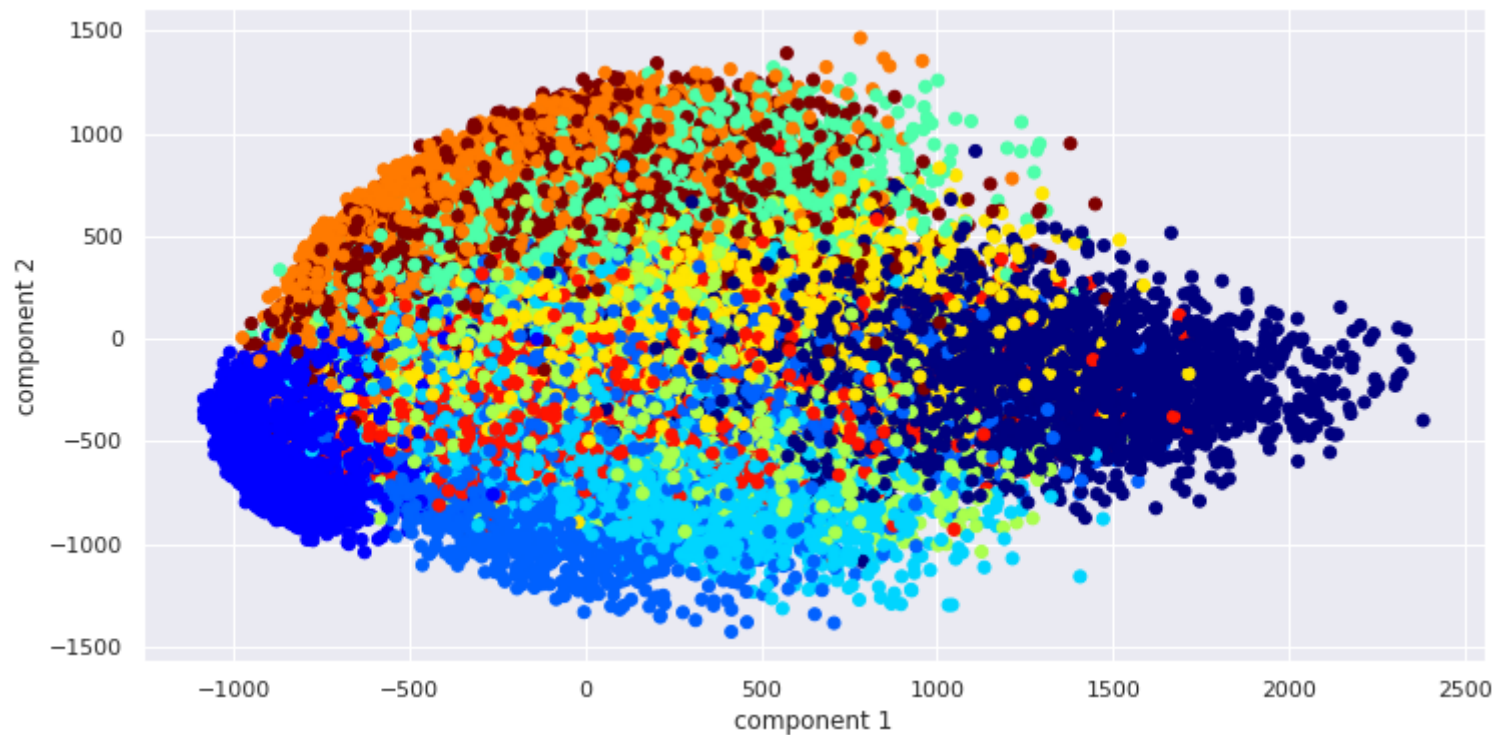
Suppose we retained only 2 synthetic features.

We would be able to plot each sample in two dimensions using the transformed coordinates.

Although targets are not necessary for PCA, here we do have labels associated with the images.

Let's plot the samples and color them according to their target.

```
In [28]: _ = ush.mnist_plot_2D(X_mnist_train_reduced, y_mnist_train.astype(int))
```



Each color is a different digit.

You can see that the clustering is far from perfect

- but also surprisingly good considering we're using only 2 out of 784 features

Let's see how much variance is captured by only the first two synthetic features.

```
In [29]: cumvar_mnist = np.cumsum(pca_mnist.explained_variance_ratio_)
first_comp = 2
cumvar_first = cumvar_mnist[first_comp-1]

print("Cumulative variance of {d:d} PC's is {p:.2f}%, about {n:.1f} pixels".format(
    d=first_comp, p=100 * cumvar_first, n=cumvar_first * X_mnist_train.shape[1]))
```

Cumulative variance of 2 PC's is 16.90%, about 132.5 pixels

Is 17% good ? You bet !

With 784 original features (pixels)

- if each feature had equal importance, it would explain $1/784 = .12\%$ of the variance.

So the first synthetic feature captures the variance of 132 original features

- (assuming all were of equal importance).



Digits example

`sklearn` comes with a "digits" dataset that is like a low resolution version of MNIST: 8 by 8 pixels.

VanderPlas has a good [notebook](#) ([external/PythonDataScienceHandbook/notebooks/05.09-Principal-Component-Analysis.ipynb#PCA-as-dimensionality-reduction](#)) illustrating PCA on this much smaller dataset.

We'll give a quick take on this, but I recommend the notebook.

Let's look at the digits data

```
In [30]: from sklearn.datasets import load_digits  
digits = load_digits()  
digits.data.shape
```

```
Out[30]: (1797, 64)
```

```
In [31]: vpd = unsupervised_helper.VanderPlas()  
_ = vpd.digits_plot(digits.data)
```



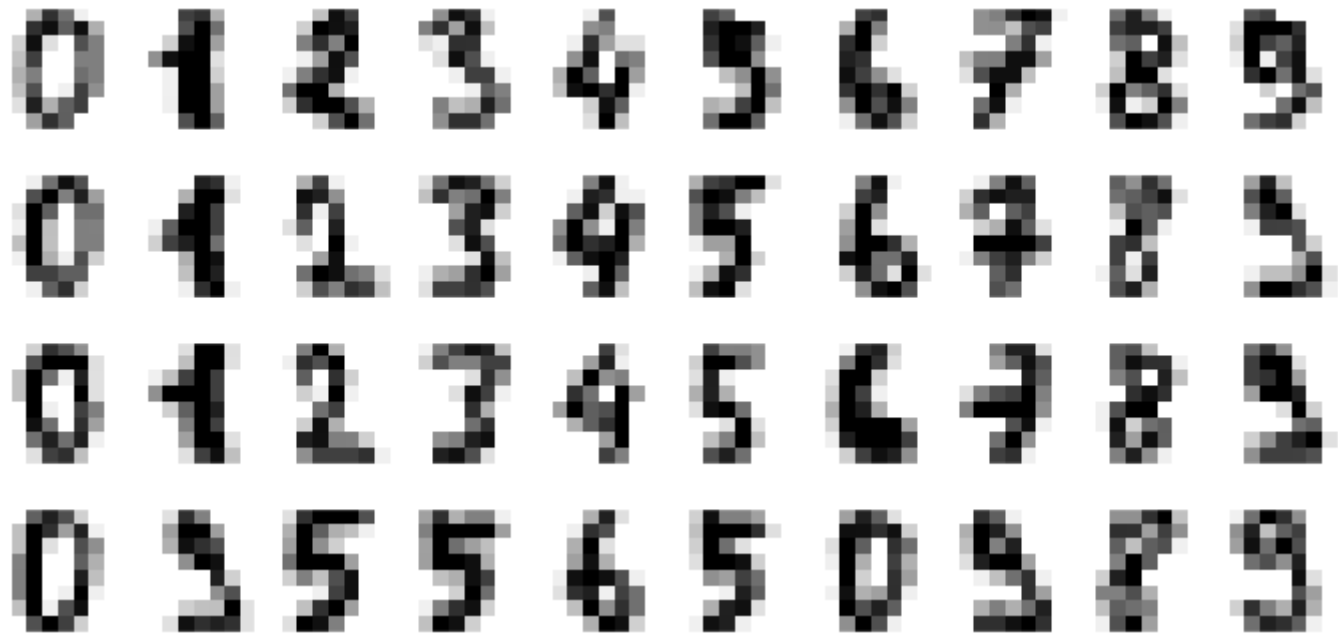
As you can see, each digit is represented by a grid of 8×8 pixels.

Let's reduce digits from 64 dimensions to only 2 and use this for clustering.

- Project to 2 synthetic features
- Use an inverse transformation to reconstruct back to original features

First, recall that if we retained all 64 sythetic features, we'd get a perfect reconstruction:

```
In [32]: digits_pca, digits_projected, digits_reconstructed = vpd.digits_reconstruction(d
digits.data)
```



Here's the reconstruction when retaining only 2 synthetic features.


```
In [33]: digits_pca2, digits_projected2, digits_reconstructed2 = vpd.digits_reconstruction(digits.data, n_components=2)
```



Reconstruction is kind of fuzzy but the digits are recognizable.

How much variance did we retain ?

```
In [34]: cumvar_digits = np.cumsum(digits_pca2.explained_variance_ratio_)
first_comp = 2
cumvar_first = cumvar_digits[first_comp-1]

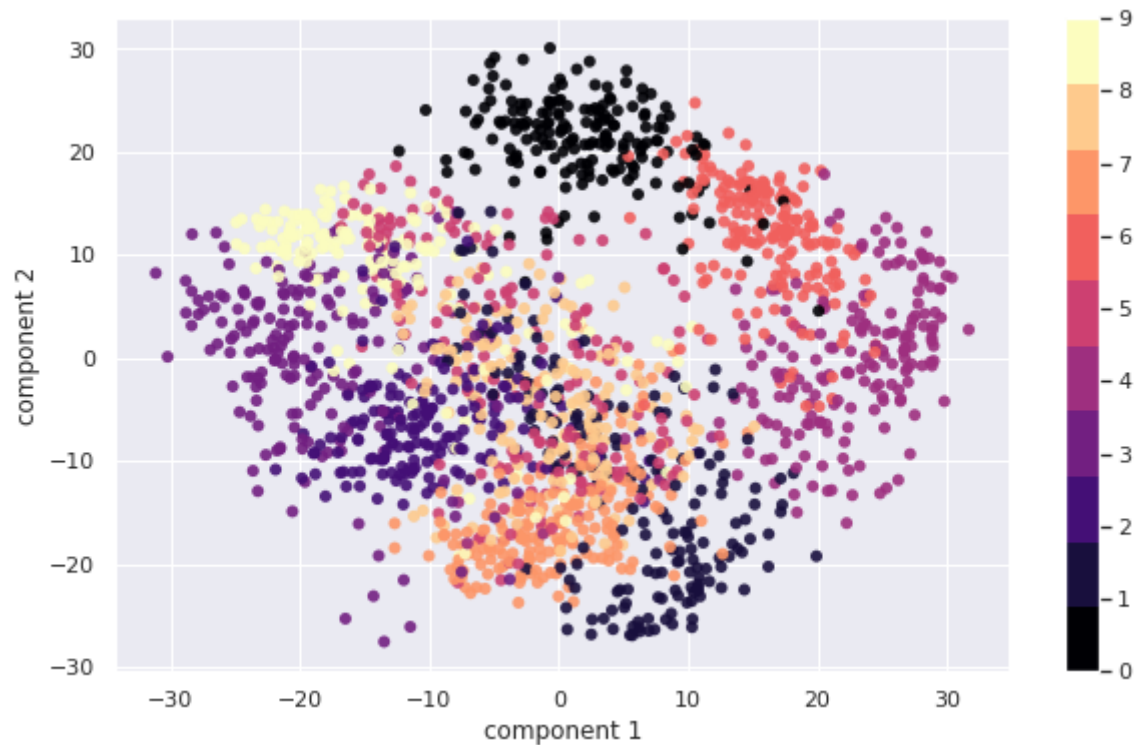
print("Cumulative variance of {d:d} PC's is {p:.2f}%, about {n:.1f} pixels".format(
    at(
        d=first_comp, p=100 * cumvar_first, n=cumvar_first * digits.data.shape[1]))
```

Cumulative variance of 2 PC's is 28.51%, about 18.2 pixels

Two synthetic features capture almost 30% of the total variance.

More importantly, let's visualize the clusters in synthetic feature (PC) space.

```
In [35]: vpd.digits_show_clustering(digits_projected2, digits.target)
```



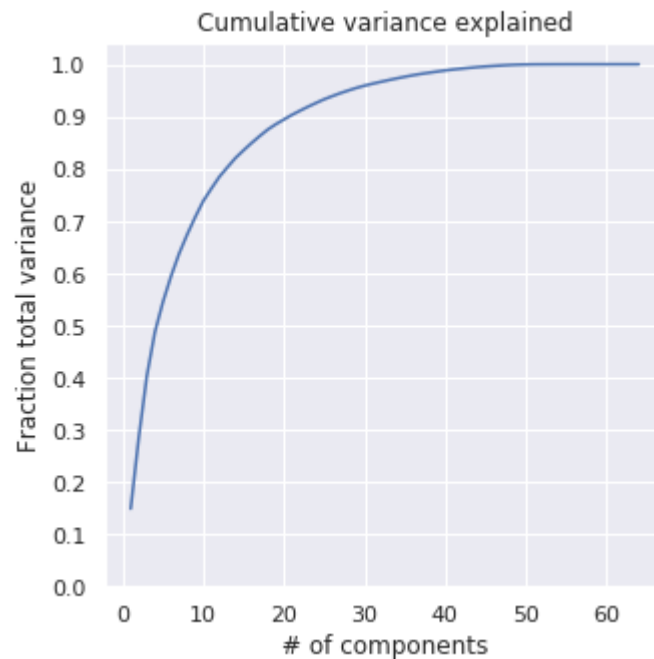
Sometimes you can interpret the PC's by looking how they separate the clusters.

I admit: don't see it here though.

What would happen if we retained more than 2 synthetic features ?

```
In [36]: _ = vpd.plot_cum_variance(digits_pca)
cumsum = np.cumsum(digits_pca.explained_variance_ratio_)
print("Cumulative variance of {nc:d} PCs is {cv:0.3f}".format(nc=2, cv=cumsum[2-1]))
```

Cumulative variance of 2 PCs is 0.285





PCA in Finance

PCA of yield curve

Litterman Scheinkman (<https://www.math.nyu.edu/faculty/avellane/Litterman1991.pdf>).

This is one of the most important papers (my opinion) in quantitative Fixed Income.

It allows us to hedge a large portfolio of bonds with a handful of instruments.

Before we show the result: why is this an important advance in Finance ?

- Imagine we had a large portfolio of bonds with many maturities.
- A common goal in Fixed Income Finance is to *immunize* (hedge) a portfolio to changes in the Yield Curve.
- A simple way to construct the hedge is to
 - find the sensitivity of each bond in the portfolio to the $n = 14$ maturities
 - sum (over bonds in the portfolio) the individual bond sensitivities
 - minus 1 times resulting portfolio sensitivity is hedge that minimizes the portfolio's exposure to Yield Curve changes

But there are transaction costs (and complexity) with $n = 14$ bonds in the hedge portfolio.

Can we do nearly as well with $n' < n$ hedge bonds ?

That's exactly what PCA is designed for: dimensionality reduction.

- in this case, reducing the number of hedge bonds
- with minimal impact on immunization goal

Let's get the history of Yield Curves and look at first few samples.

In [37]: ych = unsupervised_helper.YieldCurve_PCA()

```
# Get the yield curve data  
data_yc = ych.create_data()  
data_yc.head()
```

Out[37]:

	1M	2M	3M	6M	1J	2J	3J	4J	5J	6J	7J	8J	9J	10J
1992-02-29	0.0961	0.09610	0.0961	0.0958	0.0898	0.0864	0.0849	0.0837	0.0826	0.0817	0.0810	0.0806	0.0803	0.0804
1992-03-31	0.0970	0.09700	0.0970	0.0969	0.0912	0.0889	0.0877	0.0864	0.0852	0.0841	0.0833	0.0827	0.0823	0.0823
1992-04-30	0.0975	0.09750	0.0975	0.0975	0.0920	0.0892	0.0877	0.0862	0.0848	0.0837	0.0828	0.0822	0.0817	0.0816
1992-05-31	0.0978	0.09785	0.0979	0.0979	0.0920	0.0889	0.0874	0.0860	0.0847	0.0836	0.0828	0.0821	0.0817	0.0815
1992-06-30	0.0974	0.09745	0.0975	0.0975	0.0931	0.0904	0.0889	0.0874	0.0860	0.0848	0.0839	0.0832	0.0827	0.0825

In [38]: `data_yc.shape`

Out[38]: (287, 14)

Each example (row) has 14 features: the yields for 14 maturity points on a given date.

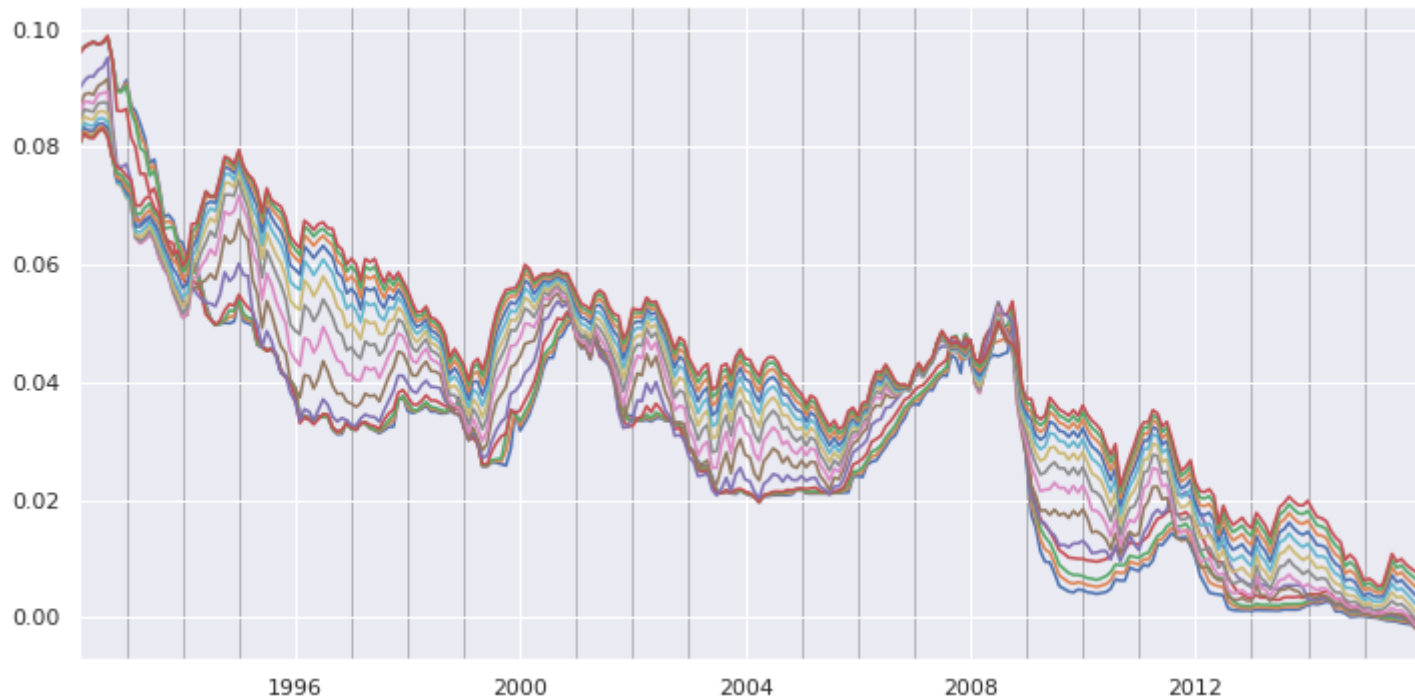
Let's plot the history of yield curves


```
In [39]: ych.plot_YC(data_yc)
```

```
/home/kjp/anaconda3/lib/python3.7/site-packages/pandas/plotting/_matplotlib/converter.py:103: FutureWarning: Using an implicitly registered datetime converter for a matplotlib plotting method. The converter was registered by pandas on import. Future versions of pandas will require you to explicitly register matplotlib converters.
```

To register the converters:

```
>>> from pandas.plotting import register_matplotlib_converters
>>> register_matplotlib_converters()
warnings.warn(msg, FutureWarning)
```



Let's perform PCA on the **changes** in Yield Curve

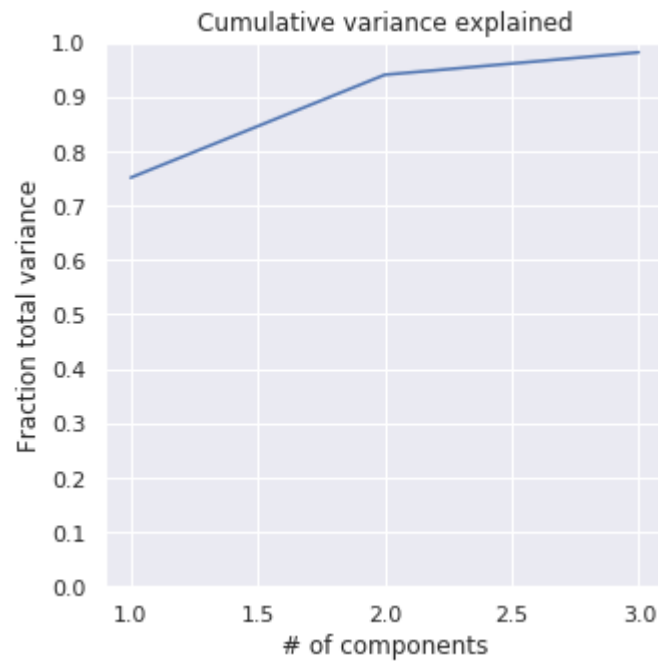
- Just like in Supervised Learning, we sometimes need to transform the data before fitting a model
- The features we feed to PCA are *yield changes* rather than yields

So $n = 14$ maturities, for m samples (many years of daily data)

How many bonds (i.e, what is the n') is "good enough" ?

The plot of cumulative variance explained, versus n' will give us an answer.

```
In [40]: pca_yc, df_pca_yc = ych.doPCA(data_yc, doDiff=True)  
_ = ych.plot_cum_variance(pca_yc)
```



Wow !

Only $n' = 3$ synthetic features capture almost all the variance of the original $n = 14$ features !

It gets even better !

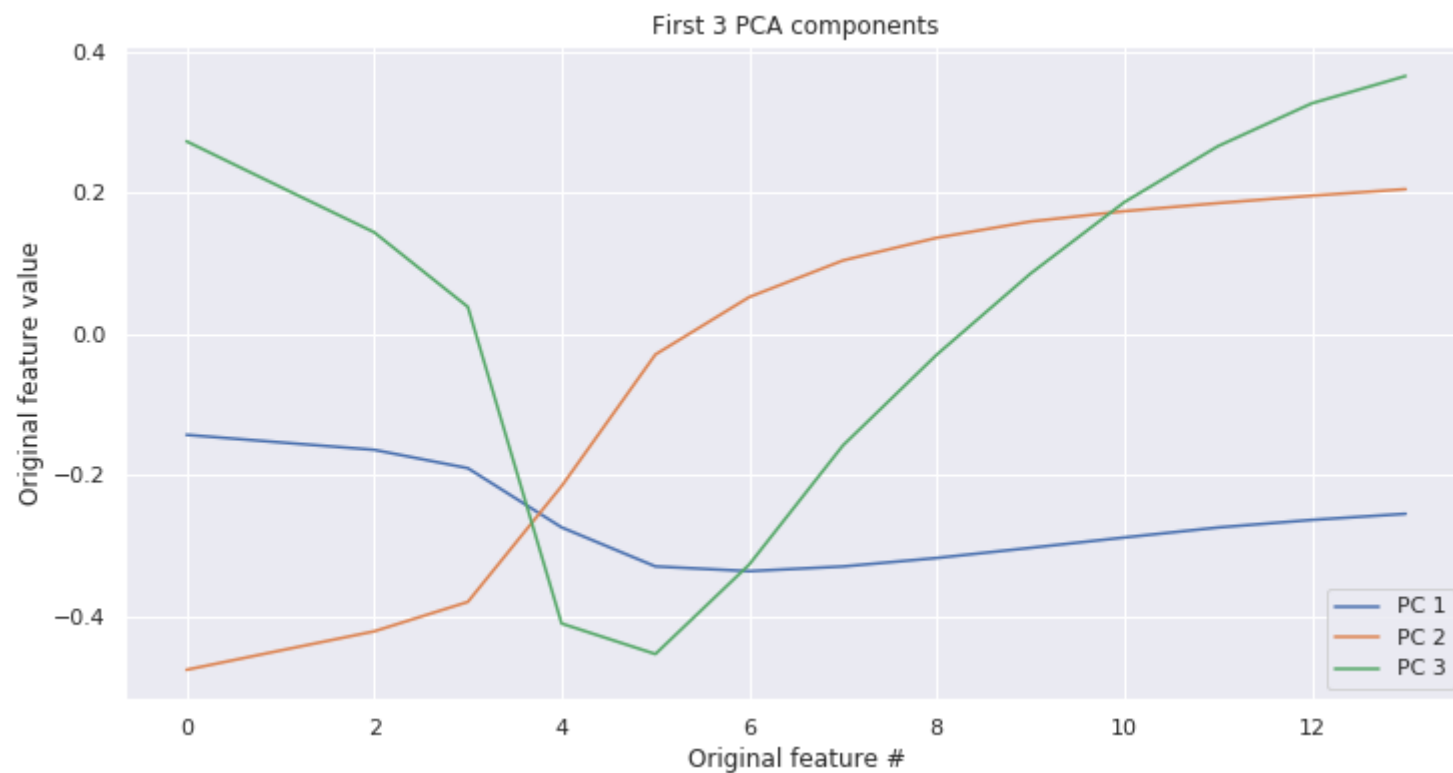
By examining the composition of the synthetic features, we can *interpret* what they are

Let's examine the effect of a 1 standard deviation move in each synthetic feature on original features.

We arrange the original features by maturity on the horizontal axis and plot the effect on the vertical.

The choice to arrange the horizontal axis by maturity is *very deliberate*, as we will see.

```
In [41]: ych.plot_components(pca_yc)
```



This is a very typical pattern in Finance:

- the first synthetic feature affects all original features with roughly equal effect
- higher synthetic features often express a *dichotomy*
 - positive effect on some original features
 - negative effect on other original features

In our case the original features are yield changes.

A unit standard deviation value of synthetic feature j (PC j)

- ($j = 1$): affects all original features (yield change) roughly equally
 - corresponds to a parallel shift in the Yield Curve
- ($j = 2$): shows a dichotomy (of yield changes) between near and far maturities
 - corresponds to the slope of the Yield Curve changing
- ($j = 3$): shows a dichotomy of yield changes of mid maturities versus extreme maturities
 - corresponds to a twist in the Yield Curve about the 5 year maturity

Recall

The synthetic features are standardized, hence 1 unit is one standard deviation.

The Σ matrix re-scales from standard deviation to original feature space.

- So can't compare the levels on the vertical axis between synthetic features.
 - Since $\sigma_1 > \sigma_2$
 - the absolute effect of synthetic feature 1 is greater than that of synthetic feature 2
 - for a 1 standard deviation move in each.

We started with a complicated sample (Yield Curve with $n = 14$ features)

- can almost completely explain (changes in the Yield Curve) with 3 intuitive market changes
 - Parallel Shift up/down of all maturities
 - Long end versus short end changes (slope)
 - Twist at intermediate maturity

That's interesting from a Machine Learning perspective but important for Finance:

- It shows how to construct efficient Hedge Portfolios.

Suppose we have a Target Portfolio consisting of long positions in many bonds, with many maturities.

Our goal is to construct another portfolio (the Hedge Portfolio)

- consisting of a small number of bonds
- with the *same* sensitivity to Yield Curve changes as the Target Portfolio

By combining a long position in the Target Portfolio with a short position in the Hedge Portfolio

- the resulting Net Portfolio is immunized (hedged) to changes in the Yield Curve

If we can re-express our Target Portfolio in terms of synthetic features

- We can construct the Hedging Portfolio as one consisting only of 3 synthetic bonds
- In quantities that mimic the sensitivity of the Target Portfolio to the first 3 PCs.
 - since the PCs are *independent*, this is easy
 - normally: you would need to solve
 - a system of 3 equations in 3 unknowns (size of position in each synthetic bond)

The catch is that each synthetic bond (feature) is a linear combination of $n = 14$ real bonds.

So we would need n real bonds to create each synthetic.

One way to deal with this is to project the synthetic features onto $n'' < 14$ real bonds.

So we come up with "approximate" versions of the synthetic features themselves.

The interpretation of Yield Curve changes guides us into a simpler hedge

- Hedge the parallel shift with the 10 year Treasury (most liquid bond in the US)
 - works because all yield changes roughly the same for first PC
- Hedge the slope with a long/short portfolio of 2 year/30 year (or 10 year) bonds
 - also liquid instruments, that are "close enough" to PC2
- Hedge the twist with a butterfly of long 2 year/10 year, short 5 year

So our Hedging Portfolio will consist of just a handful of very liquid bonds

- as opposed to n bonds, many of which may be illiquid

The primary way of hedging Fixed Income portfolios prior to this was by *duration* hedging

- assuming that yields across all maturities moved the same amount
- using a single liquid bond as the Hedge Portfolio

The PCA verified that the simple, intuitive hedge was actually the most important hedge to make !

Finance details

This section has little to do with Machine Learning but quite a bit to do with Fixed Income Finance.

- We have captured changes in *yield* of a bond
- To hedge *price* changes (our goal) we still need to translate a yield change to a price change

- For bond b with price P_b and yield y_b
 - we need $\frac{\partial P_b}{\partial y_b}$
 - the change in Price of bond b per unit change in its yield
 - this is known as the bond's *duration*
 - if we hold $\#_b$ units of bond b in a Portfolio
 - the bond's contribution to portfolio price change is $\#_b$ times the above sensitivity
 - sometimes more convenient to compute the *percent price* change per unit yield change
 - size of the hedge now in *number of dollars* rather than *number of bonds*

PCA of the SP 500

The same analysis that we did for the Bond Universe works for other instruments.

Consider a universe of all stocks in a particular stock universe.

We can perform PCA on the returns (percent changes) of each stock

- discover the common factors affecting all stocks in the universe.
- $n = 500$ features, for each example (one day return)

We don't have time to do it here, but the first components of many universes tends to be

- a first component that has roughly equal impact
 - PC1 is almost an equally weighted portfolio of all stocks
- higher components expressing dichotomies
 - cyclical stocks versus non-cyclical
 - large cap versus small cap
 - industry versus other industry

Interpreting the PC's

The key in our interpretation of the PCs for the Yield Curve

- choice of ordering our original features by sorted maturity.

Had we chosen some other arrangement of the horizontal axis, we may not have seen the pattern.

So how do we find the "right" pattern ?

- Assign each original feature a set of attributes
 - bonds: maturity
 - stocks: industry, market capitalization
- Propose a "theory" about how the value of an original feature will respond to a level of the PCj
 - The theory should relate the attribute of an original feature to the value of feature
- As a horizontal axis: sort by the attribute proposed by your theory
 - Stocks:
 - ($j = 1$): arbitrary arrangement works, since all stocks responds equally
 - ($j > 1$): cluster stocks by attribute
 - sort by market capitalization
 - group by industry: first all Industrials, then all Techs, etc.

Let \hat{u} be the $1 \times n$ vector of all 0's except for a 1 at position i

$$\begin{aligned}\hat{u}_i &= 1 \\ \hat{u}_j &= 0, \quad j \neq i\end{aligned}$$

Then

$$s = \hat{u}IV^T$$

is a $1 \times n$ vector whose elements are the effect of a one standard deviation shock in synthetic factor i on each original feature.

$s_j = \text{change in } X_j \text{ for a 1 standard deviation change in } \tilde{X}_i$

$$s' = \hat{u}\Sigma V^T$$

is a $1 \times n$ vector that is scaled by the actual standard deviation of \tilde{X}_i since the absolute size of a 1 standard deviation change in \tilde{X}_i is σ_i .

That is: row i of V^T is the effect on each original feature of a one standard deviation shock in synthetic feature \tilde{X}_i .



Recommender Systems: Pseudo SVD

There is another interesting use of Matrix Factorization that we will briefly review.

It will show both a case study and interesting extension of SVD.

Netflix Prize competition

- Predict user ratings for movies
- Dataset
 - Ratings assigned by users to movies: 1 to 5 stars
 - 480K users, 18K movies; 100MM ratings total
- \$1MM prize
- awarded to team that beat Netflix existing prediction system by at least 10 percentage points

User preference matrix

We will try to use same language as PCA (examples, features, synthetic features)

- but map them to Netflix terms
 - Examples: Viewers
 - Features: Movies ("items")

Matrix \mathbf{X} : user rating of movies

$X_j^{(i)}$ is i^{th} user's rating of movie j

X is huge: $m * n$

- $m = .5$ million viewers
- $n = 18,000$ items (movies).

About 9 billion entries for a full matrix !

Idea: Linking Viewer to Movies via concepts

- Come up with your own "concepts" (synthetic features)
 - concept = attribute of a movie
 - map user preference to concept
 - map movie style to concept
 - supply and demand:
 - user demands concept, movie provides concept

Human defined concepts

- style: action, adventure, comedy, sci-fi
- actor
- typical audience segment

- Create user profile P : maps user to concept
- Create item profile Q : maps movies (features, items) to concept
- $\mathbf{X} = PQ^T$

One advantage of the $\mathbf{X} = PQ^T$ approach is a big space reduction.

With $k \leq n$ concepts:

- \mathbf{X} is $m \times n$
- P is $m \times k$
- Q is $n \times k$

So idea is to factor X and discover $k \leq n$ synthetic features (concepts, "latent factors")

- $\mathbf{X} = PQ^T$
 - fit the model on training examples, i.e., perform decomposition
- Given a new user (test example), predict rating for an unseen movie:
 - Get a partial vector of user's ratings (defined for movies seen by user)
 - original features
 - Project onto concepts (synthetic features):
 - Inverse transform to get back complete vector in original feature space

SVD to discover concept

Why let a human guess what ML can discover ?

Use SVD to discover the k "best" synthetic features, rather than leaving it to a person.

Factor \mathbf{X} by SVD !

$$\mathbf{X} = U\Sigma * V^T$$

Let $P = U\Sigma, Q = V$

Anyone spot the problem(s) ?

- First: a matrix with 9 billion entries is a handful !
- As you can imagine, any single user views only a fraction of the m movies
 - \mathbf{X} is very sparse
 - Of the 9 billion potential entries in the full matrix \mathbf{X} , we only have 100 million defined.

How can we perform SVD on matrix with missing values ?

- Missing value imputation
 - not attractive: most values are missing
- Pseudo SVD

The ML mantra

- find a (cool ?) cost function that describes a solution to your problem
- Use Gradient Descent to solve

Pseudo SVD Cost function

The Frobenius Norm (used above) modified to exclude missing values

$$Cost(\mathbf{X}', \mathbf{X}) = \sum_{\substack{1 \leq i \leq m, \\ 1 \leq j \leq n \\ \mathbf{X}_j^{(i)} \text{ defined}}} \left(\mathbf{X}_j^{(i)} - \mathbf{X}'_j^{(i)} \right)^2$$

where $\mathbf{X}' = PQ^T$

Interpret as Reconstruction Error

Pseudo SVD algorithm

- Define $\mathbf{X}' = PQ^T$
- Take analytic derivatives of $Cost(X', X)$ with respect to
 - $P_j^{(i)}$ for $1 \leq i \leq m, 1 \leq j \leq k$
 - $Q_j^{(i)}$ for $1 \leq i \leq m, 1 \leq j \leq k$
- Initialize elements of P, Q randomly.
- Use Gradient Descent to solve for optimal entries of P, Q .
 - find entries of P, Q such that product matches non-empty part of \mathbf{X}

Note

- No guarantee that the P, Q obtained are
 - orthonormal, etc. (which SVD would give you)

But SVD won't work for \mathbf{X} with missing values.

Filling in missing values

Once you have P, Q

- to predict a missing rating for user i movie j :

$$\hat{r}_{j,i} = q^{(i)} \cdot p_j^T$$

- $q^{(i)}$ is row i of Q
- p_j is column j of P^T

Some intuition

The rating vector of a user may have missing entries.

But we can still project to synthetic feature space based on the non-empty entries.

The projection winds up in a "neighborhood" of concepts.

Inverse transformation

- gets us to a completely non-empty rating vector that is a resident of this neighborhood.

Example

User rates

- Sci-Fi movies A and B very highly
- Does not rate Sci-Fi movie C.

Since A,B, C express same concept (Sci-Fi) they will be close in synthetic feature space.

Hence, the implied rating of User for movie C will be close to what other users rate C.

In [42]: `print("Done")`

Done