

# CPSC 340: Machine Learning and Data Mining

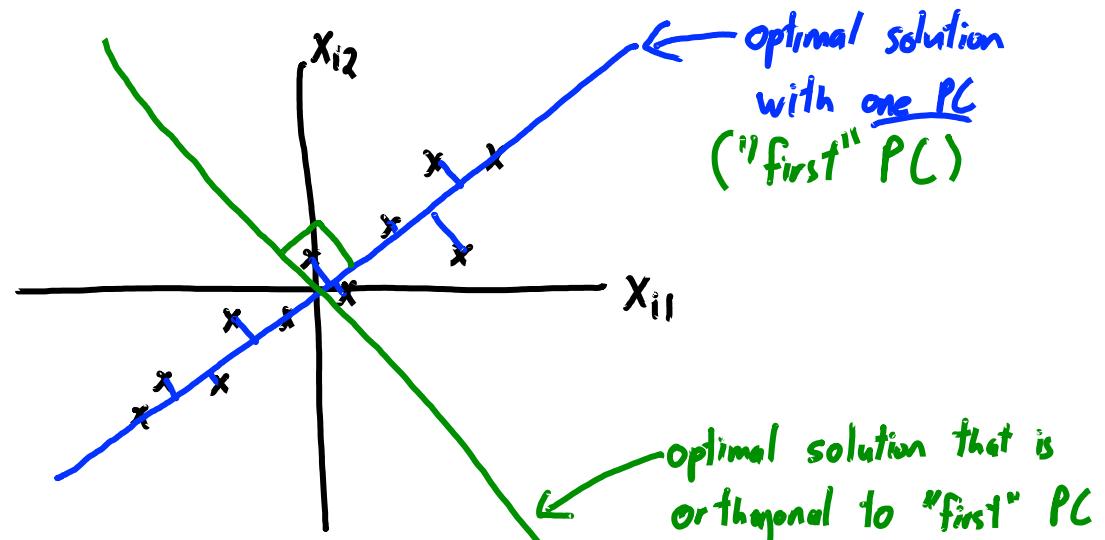
## Sparse Matrix Factorization

# Admin

- Assignment 5:
  - Due next Friday.

# Last Time: PCA with Orthogonal/Sequential Basis

- When  $k = 1$ , PCA has a **scaling problem**.
- When  $k > 1$ , have **scaling, rotation, and label switching**.
  - Standard fix: use **normalized orthogonal rows  $W_c$**  of ' $W$ '.
$$\|w_c\|=1 \quad \text{and} \quad w_c^T w_{c'} = 0 \quad \text{for } c' \neq c$$
  - And **fit the rows in order**:
    - First row "explains the most variance" or "reduces error the most".



# Application: Face Detection

- Consider problem of face detection
- Classic methods use “eigenfaces” as basis:
  - PCA applied to images of faces.

# Eigenfaces

- Collect a bunch of images of faces under different conditions:



Each row of  $X$  will be pixels in one image:

$$X =$$

If have ' $n$ ' images that are ' $m$ ' by ' $m$ ' then  $X$  is ' $n$ ' by  $m^2$ .

# Eigenfaces

Compute mean  $\mu_j$  of each column.



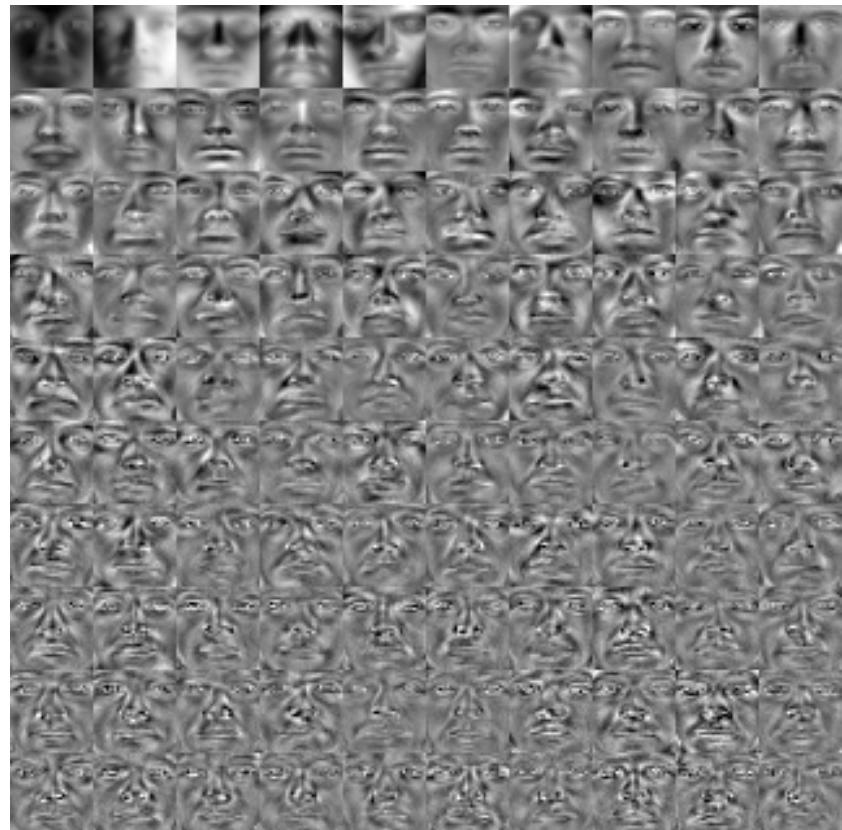
Each row of  $X$  will be pixels in one image:

$$X = \begin{bmatrix} x_1 - \mu \\ x_2 - \mu \\ \vdots \\ x_n - \mu \end{bmatrix}$$

Replace each  $x_{ij}$  by  $x_{ij} - \mu_j$

# Eigenfaces

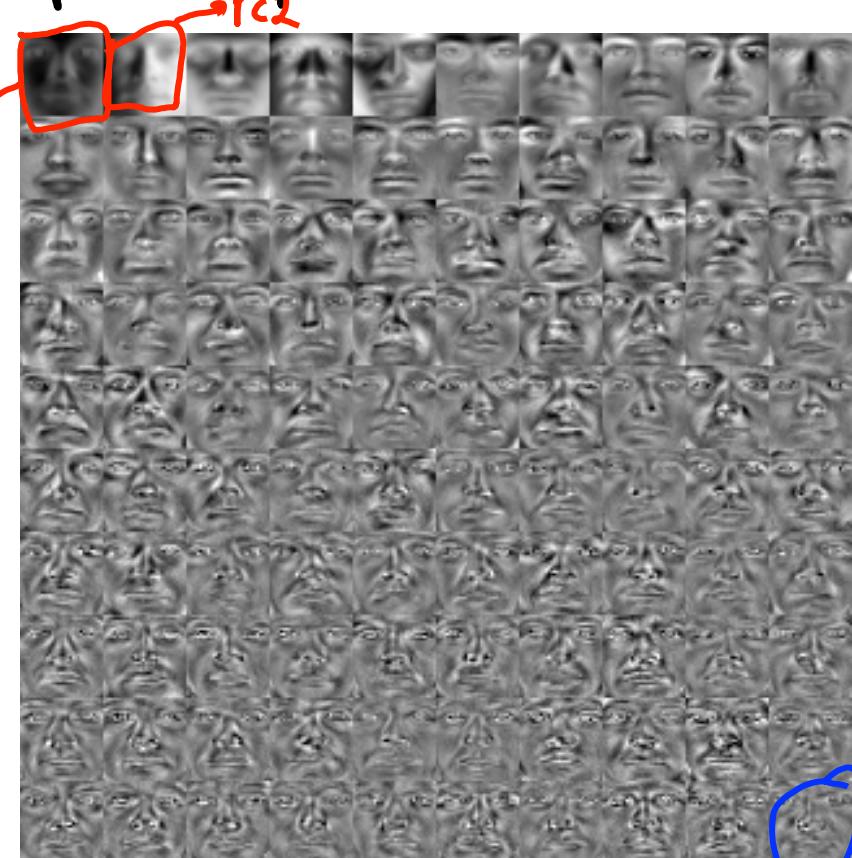
Compute top 'k' PCs on centered data: Each row of  $X$  will be pixels in one image:



$$X = \begin{bmatrix} x_1 - \mu \\ x_2 - \mu \\ \vdots \\ x_n - \mu \end{bmatrix}$$

# Eigenfaces

Compute top 'k' PCs on centered data:



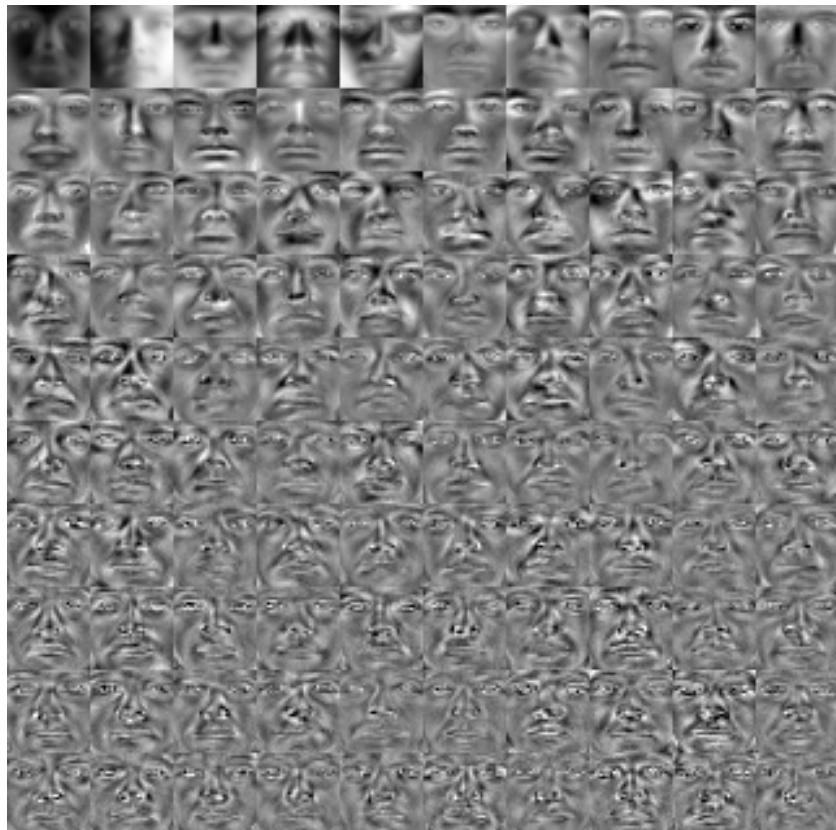
Note that these are "signed" images.



"gray" represents values close to 0.  
"dark" represents negative values  
"bright" represents positive values

# Eigenfaces

Compute top 'k' PCs on centered data:



"Eigenface" representation:

$$\hat{x}_i = \mu + z_{i1} \text{PC1} + z_{i2} \text{PC2} + z_{i3} \text{PC3} + \dots$$

$\hat{x}_i$        $\mu$        $z_{i1}$        $z_{i2}$        $z_{i3}$        $\dots$

PC1  
(first row of  $W$ )

# Eigenfaces

100 of the original faces:



"Eigenface" representation:

$$\hat{x}_i = \mu + z_{i1} \text{PC1} + z_{i2} \text{PC2} + z_{i3} \text{PC3} + \dots$$

$\hat{x}_i$        $\mu$        $z_{i1}$        $\text{PC1}$   
 $z_{i2}$        $\text{PC2}$   
 $z_{i3}$        $\text{PC3}$   
 $\vdots$        $\vdots$

(first row of  $W$ )

# Eigenfaces

Reconstruction with  $k=0$



Variance explained: 0%

“Eigenface” representation:

$$\hat{x}_i = \mu + z_{i1} \text{PC1} + z_{i2} \text{PC2} + z_{i3} \text{PC3} + \dots$$

$\hat{x}_i$        $\mu$

$z_{i1}$        $z_{i2}$        $z_{i3}$        $\dots$

PC1      PC2      PC3

(first row of  $W$ )

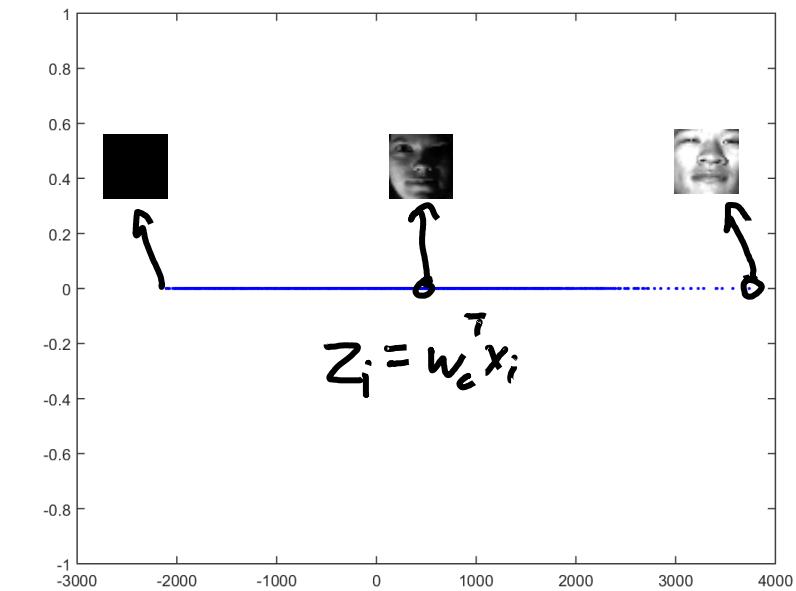
# Eigenfaces

Reconstruction with  $k=1$



Variance explained: 34%

PCA Visualization:



"Eigenface" representation:

$$\hat{x}_i = \mu + z_{i1} \text{PC1} + z_{i2} \text{PC2} + z_{i3} \text{PC3} + \dots$$

The term  $\hat{x}_i$  is circled in red. The term  $\mu$  is circled in red. The term  $z_{i1} \text{PC1}$  is circled in red. A red bracket under the terms  $\text{PC1}, \text{PC2}, \text{PC3}, \dots$  is labeled "first row of  $W$ ".

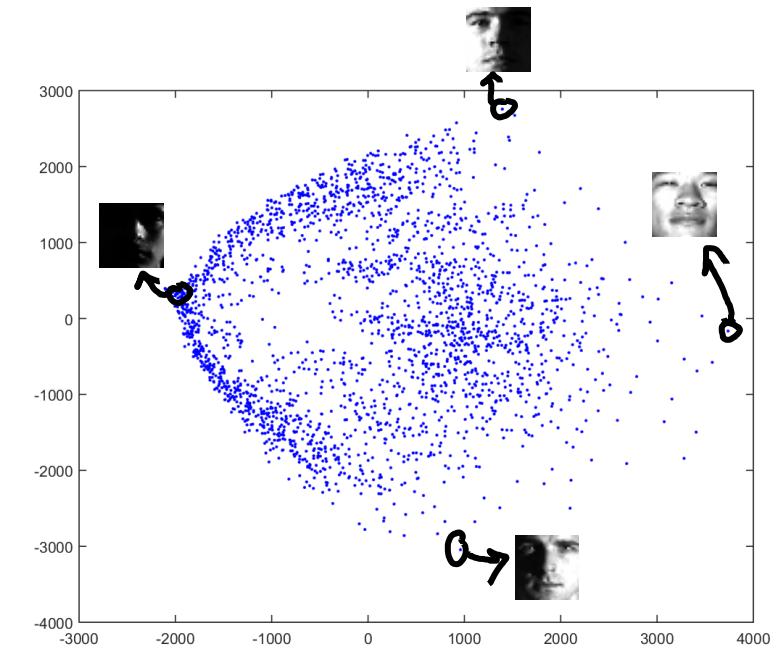
# Eigenfaces

Reconstruction with  $k=2$



Variance explained: 71%

PCA Visualization:



"Eigenface" representation:

$$\hat{x}_i = \mu + z_{i1} \text{PC1} + z_{i2} \text{PC2} + z_{i3} \text{PC3} + \dots$$

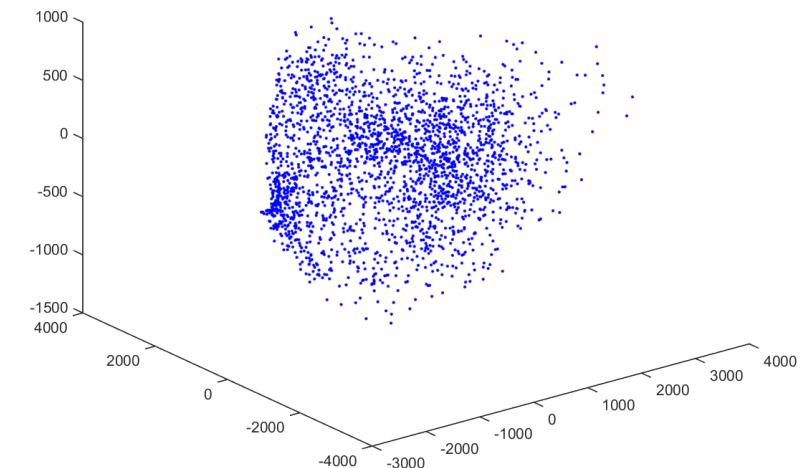
# Eigenfaces

Reconstruction with  $k=3$



Variance explained: 76%

PCA Visualization:



"Eigenface" representation:

$$\hat{x}_i = \mu + z_{i1} \text{PC1} + z_{i2} \text{PC2} + z_{i3} \text{PC3} + \dots$$

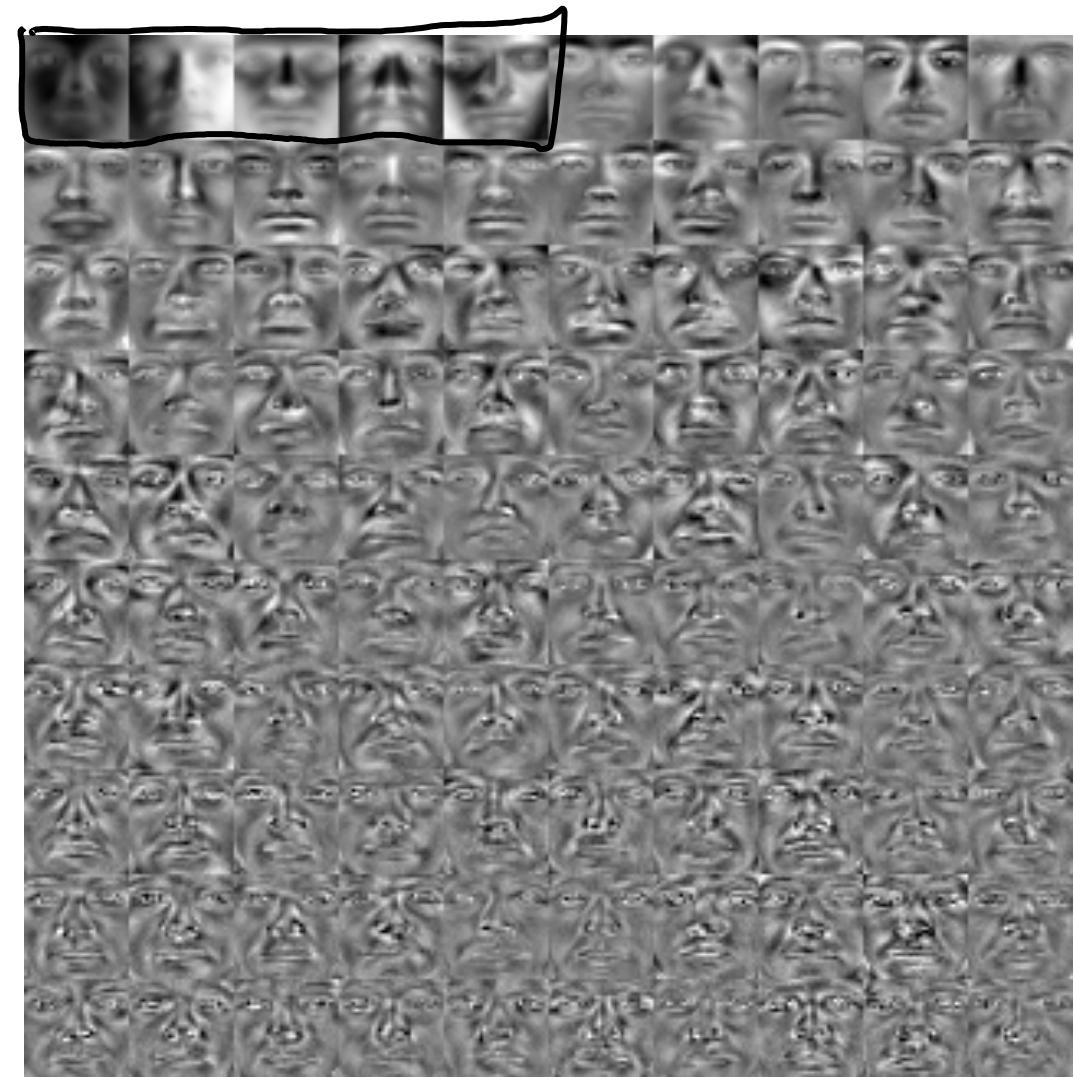
The equation shows the "Eigenface" representation of a face  $\hat{x}_i$ . It consists of a mean face  $\mu$  plus a linear combination of three principal components (eigenfaces)  $\text{PC1}$ ,  $\text{PC2}$ , and  $\text{PC3}$ , plus a residual term represented by ellipses. A red bracket groups the terms  $\mu + z_{i1} \text{PC1} + z_{i2} \text{PC2} + z_{i3} \text{PC3}$ . Below this bracket, another red bracket groups the term  $(\text{first row of } w_i)$ .

# Eigenfaces

Reconstruction with  $k=5$



Variance explained: 80%

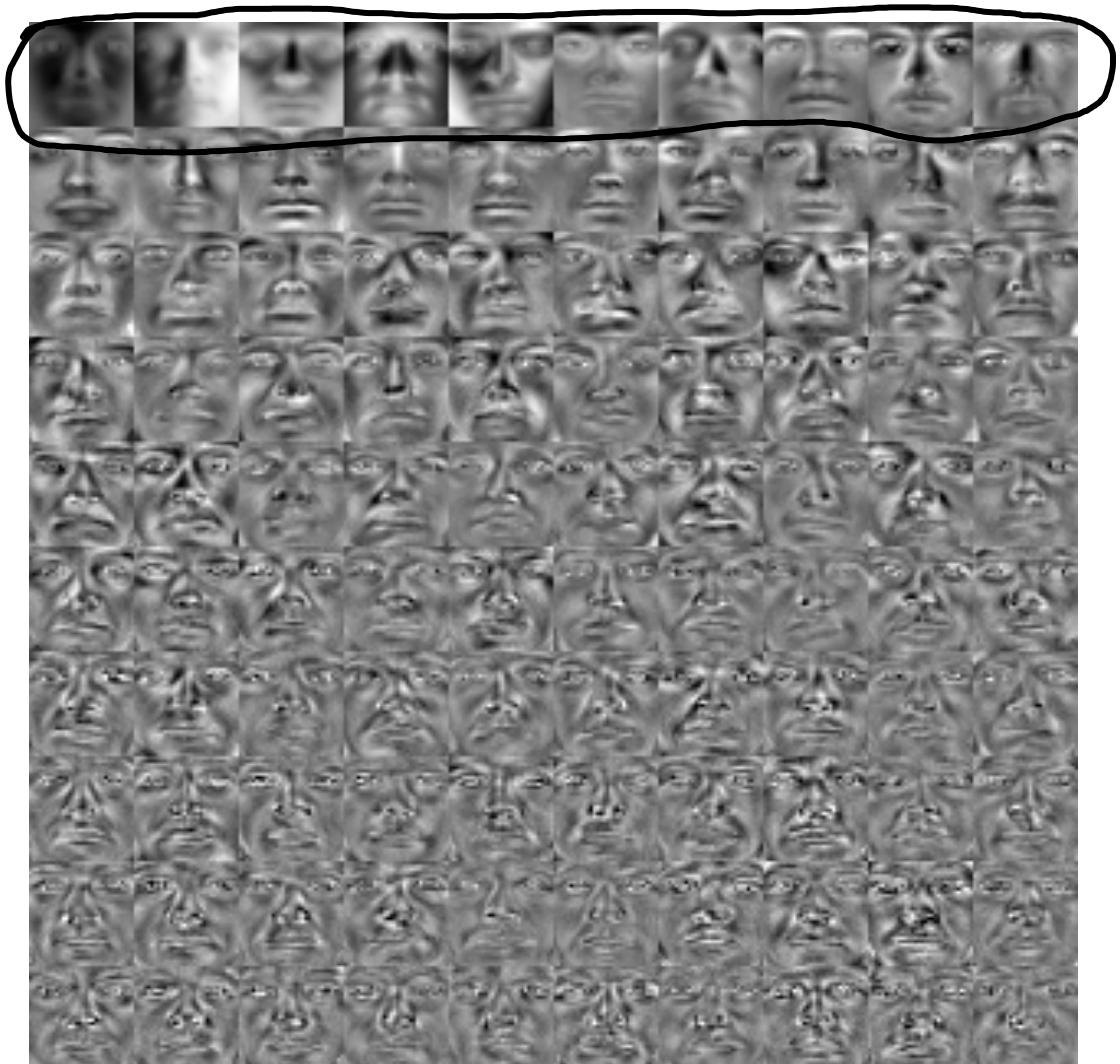


# Eigenfaces

Reconstruction with  $k=10$



Variance explained: 85%

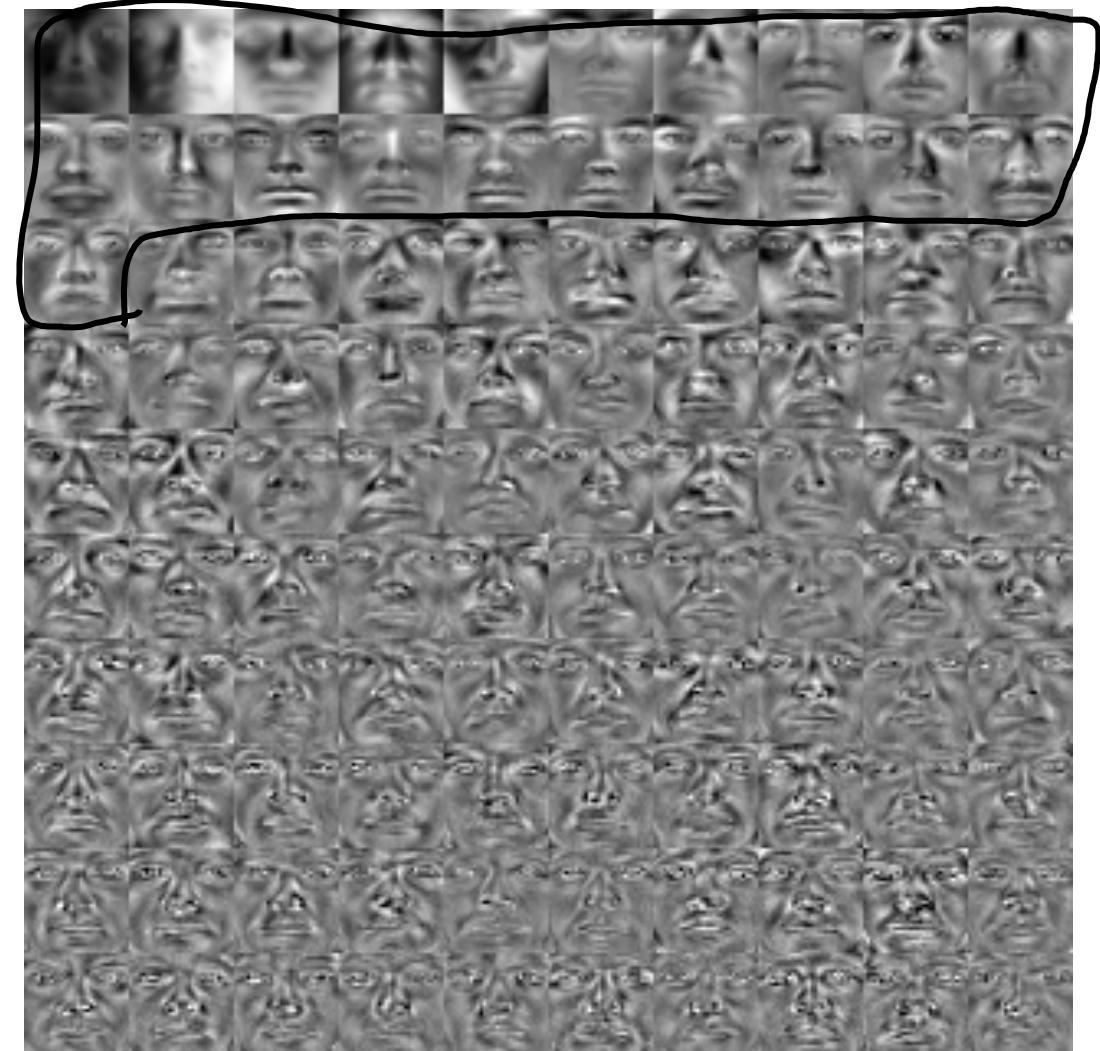


# Eigenfaces

Reconstruction with  $k=21$



Variance explained: 90%

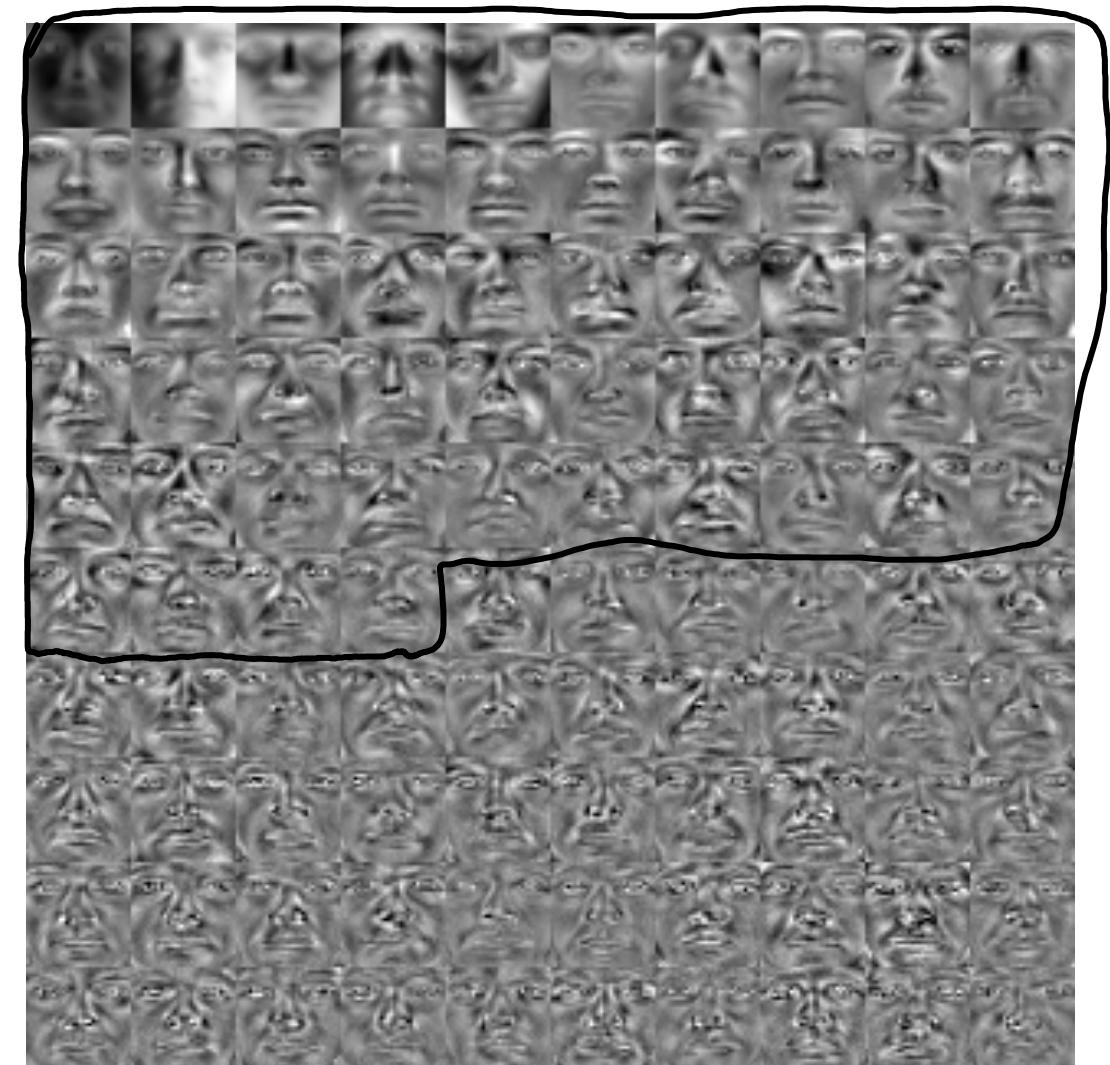


# Eigenfaces

Reconstruction with  $k=54$



Variance explained: 95%

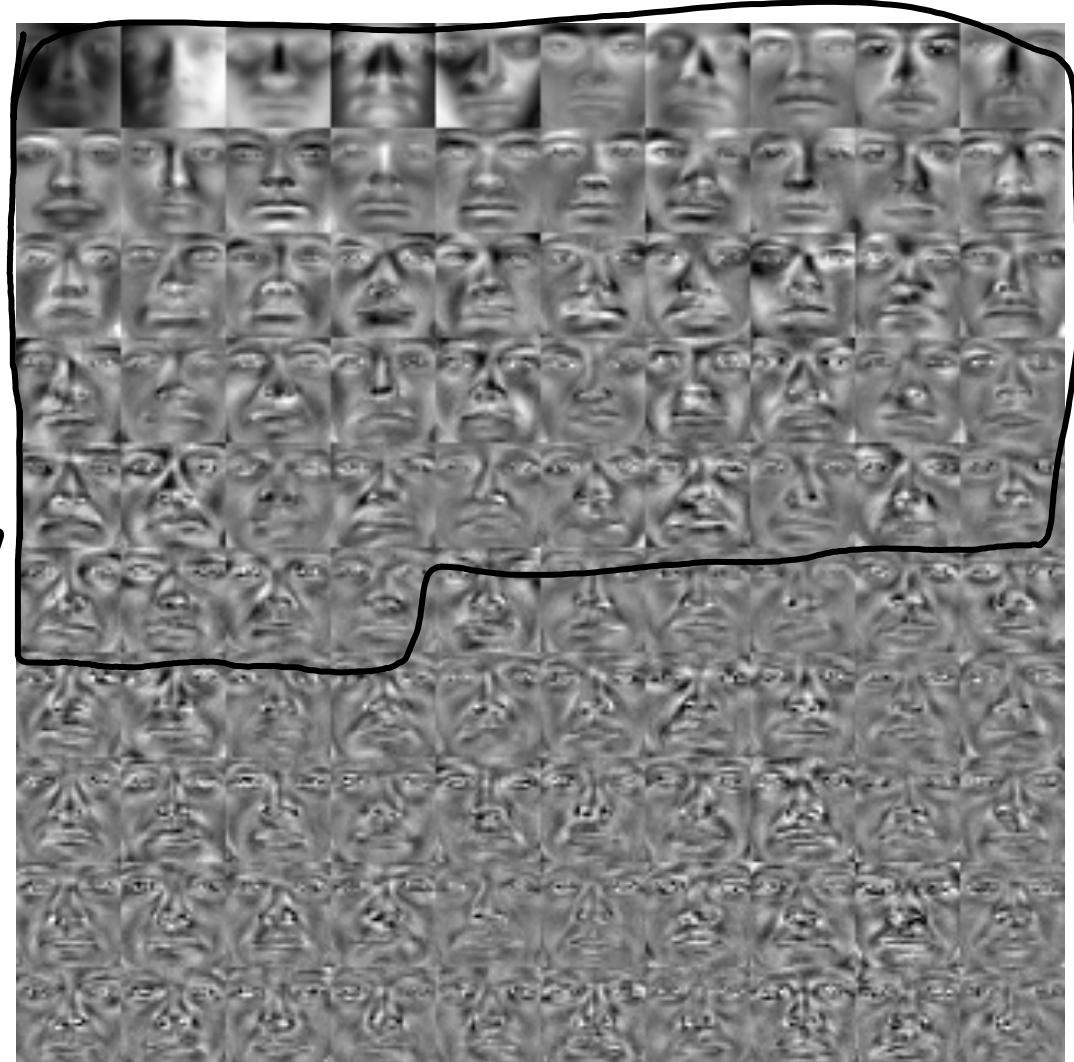


# Eigenfaces

Original Images again:



Plus these  
"eigenfaces"  
and  
the  
mean.



We can replace 1024  $x_i$  values by 54  $z_i$  values

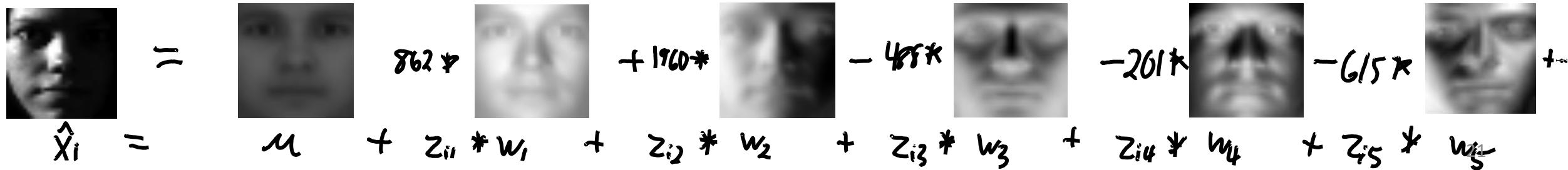
# VQ vs. PCA vs. NMF

- But how *should* we represent faces?
  - Vector quantization (k-means).
    - Replace face by the **average face** in a cluster.
    - **Can't distinguish between people** in the same cluster (only 'k' possible faces).

$$\hat{x}_i = 0 * \text{[face 1]} + 0 * \text{[face 2]} + 0 * \text{[face 3]} + 1 * \text{[face 4]} + 0 * \text{[face 5]} + 0 * \text{[face 6]} + \dots$$
$$\hat{x}_i = z_{i1} * w_1 + z_{i2} * w_2 + z_{i3} * w_3 + z_{i4} * w_4 + z_{i5} * w_5 + z_{i6} * w_6 + \dots$$

# VQ vs. PCA vs. NMF

- But how *should* we represent faces?
  - Vector quantization (k-means).
  - PCA (orthogonal basis).
    - Global average plus linear combination of “eigenfaces”.
    - But “eigenfaces” are **not intuitive** ingredients for faces.
      - PCA tends to use positive/negative **cancelling** bases.

$$\hat{x}_i = \mu + z_{i1} * w_1 + z_{i2} * w_2 + z_{i3} * w_3 + z_{i4} * w_4 + z_{i5} * w_5 + \dots$$


# VQ vs. PCA vs. NMF

- But how *should* we represent faces?
  - Vector quantization (k-means).
  - PCA (orthogonal basis).
  - NMF (non-negative matrix factorization):
    - Instead of orthogonality/ordering in W, require W and Z to be non-negativity.
    - Example of “sparse coding”:
      - The  $z_i$  are sparse so each face is coded by a small number of neurons.
      - The  $w_c$  are sparse so neurons tend to be “parts” of the object.

$$\begin{aligned} \text{Image} &= 4,2 * \text{Image}_1 + 0 * \text{Image}_2 + 0 * \text{Image}_3 + 3,3 * \text{Image}_4 + 0 * \text{Image}_5 + \dots \\ \hat{x}_i &= z_{i1} * w_1 + z_{i2} * w_2 + z_{i3} * w_3 + z_{i4} * w_4 + z_{i5} * w_5 + \dots \end{aligned}$$

# Warm-up to NMF: Non-Negative Least Squares

- Consider our usual **least squares** problem:

$$f(w) = \frac{1}{2} \sum_{i=1}^n (w^\top x_i - y_i)^2$$

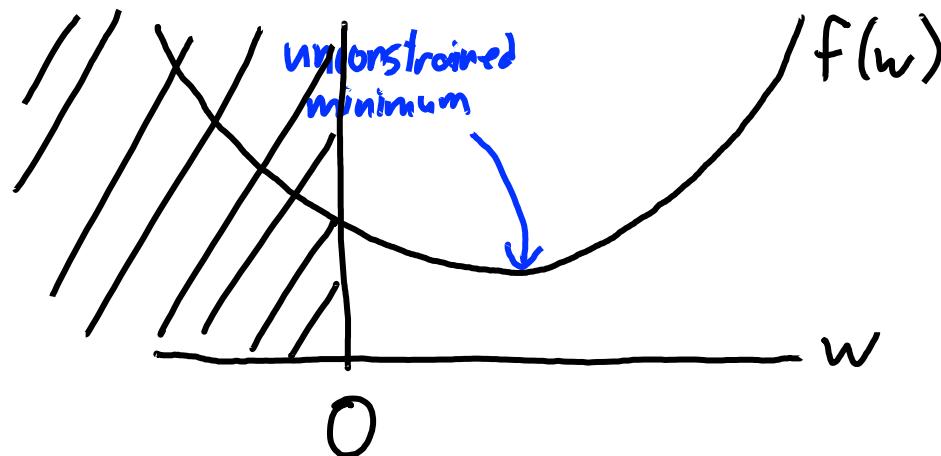
- But assume  $y_i$  and elements of  $x_i$  are **non-negative**:
  - Could be sizes ('height', 'milk', 'km') or counts ('vicodin', 'likes', 'retweets').
- Assume we want elements of ' $w$ ' to be **non-negative**, too:
  - Maybe no sensible interpretation to negative weights.
- Non-negativity leads to sparsity...

# Sparsity and Non-Negative Least Squares

- Consider 1D non-negative least squares objective:

$$f(w) = \frac{1}{2} \sum_{i=1}^n (w x_i - y_i)^2 \text{ with } w \geq 0$$

- Plotting the (constrained) objective function:



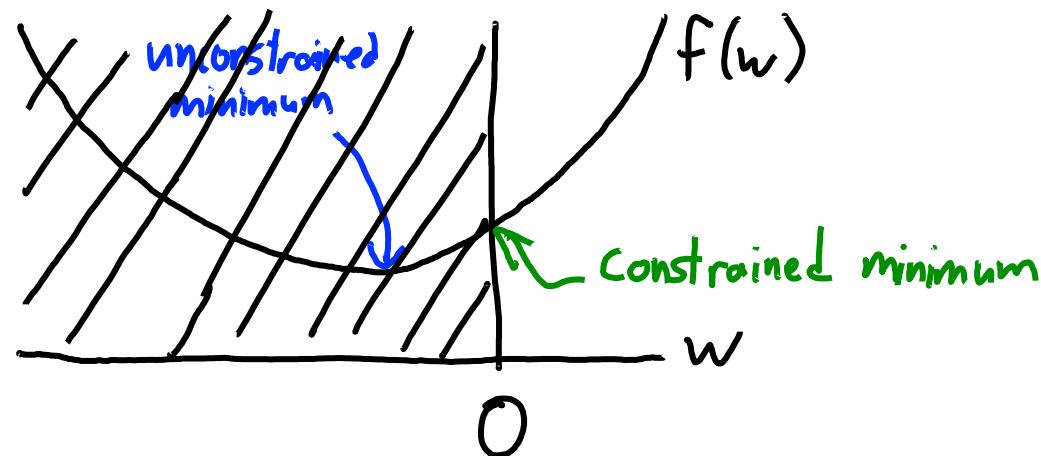
- In this case, non-negative solution is least squares solution.

# Sparsity and Non-Negative Least Squares

- Consider 1D non-negative least squares objective:

$$f(w) = \frac{1}{2} \sum_{i=1}^n (w x_i - y_i)^2 \text{ with } w \geq 0$$

- Plotting the (constrained) objective function:



- In this case, non-negative solution is  $w = 0$ .

# Sparsity and Non-Negativity

- Similar to L1-regularization, non-negativity leads to sparsity.
  - Also regularizes:  $w_j$  are smaller since can't "cancel" out negative values.
- How can we minimize  $f(w)$  with non-negative constraints?
  - Naive approach: solve least squares problem, set negative  $w_j$  to 0.  
Compute  $w = (X^\top X)^{-1} (X^\top y)$
  - Set  $w_j = \max\{0, w_j\}$
  - This is correct when  $d = 1$ .
  - Doesn't make sense when  $d \geq 2$ .
    - Consider two collinear or almost collinear features, with  $w_1=10$  and  $w_2=-10$
    - Setting  $w_1=w_2=0$  might be OK, but setting  $w_1=10$  and  $w_2=0$  is wrong.

# Sparsity and Non-Negativity

- Similar to L1-regularization, non-negativity leads to sparsity.
  - Also regularizes:  $w_j$  are smaller since can't "cancel" out negative values.
- How can we minimize  $f(w)$  with non-negative constraints?
  - A correct approach is projected gradient algorithm:
    - Run a gradient descent iteration:

$$w^{t+1} = w^t - \alpha^t \nabla f(w^t)$$

- After each step, set negative values to 0.

$$w_j^{t+1} = \max\{0, w_j^{t+1}\}$$

- Repeat.

# Sparsity and Non-Negativity

- Similar to L1-regularization, non-negativity leads to sparsity.
  - Also regularizes:  $w_j$  are smaller since can't "cancel" out negative values.
- How can we minimize  $f(w)$  with non-negative constraints?
  - A correct approach is projected gradient algorithm:

$$w^{t+1} = w^t - \alpha^t \nabla f(w^t) \quad w_j^{t+1} = \max\{0, w_j^{t+1}\}$$

- Similar properties to gradient descent:
  - Guaranteed decrease of 'f' if  $\alpha_t$  is small enough.
  - Reaches local minimum under weak assumptions (global minimum for convex 'f').
    - Least squares objective is still convex when restricted to non-negative variables.
  - Generalizations allow things like L1-regularization instead of non-negativity.

(findMinL1)

# Projected-Gradient for NMF

- Back to the non-negative matrix factorization (NMF) objective:

$$f(W, Z) = \sum_{i=1}^n \sum_{j=1}^d ((w^j)^T z_i - x_{ij})^2 \quad \text{with } w_{cj} \geq 0 \text{ and } z_{ij} \geq 0$$

- Different ways to use projected gradient:

- Alternate between projected gradient steps on 'W' and on 'Z'.
- Or run projected gradient on both at once.
- Or sample a random 'i' and 'j' and do stochastic projected gradient.

Set  $z_i^{t+1} = z_i^t - \alpha^t \nabla_{z_i} f(W, Z)$  and  $(w^j)^{t+1} = (w^j)^t - \alpha^t \nabla_{w^j} f(W, Z)$  for selected  $i$  and  $j$

- Non-convex and (unlike PCA) is sensitive to initialization.

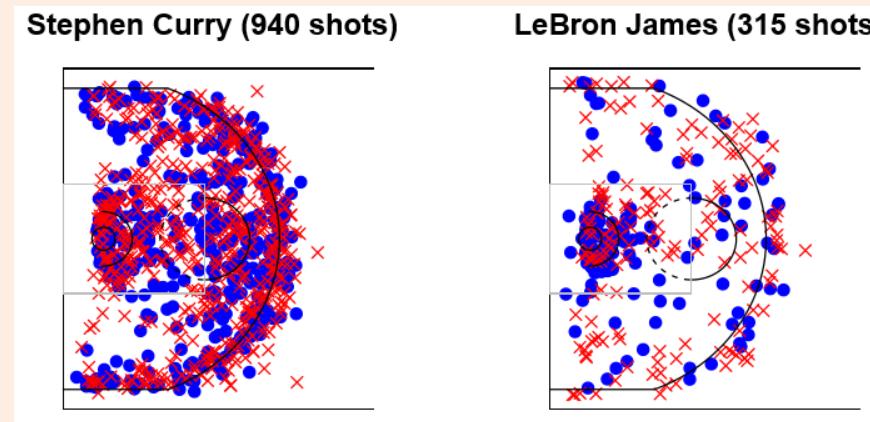
- Hard to find the global optimum.
- Typically use random initialization.

(keep other values of W and Z fixed)

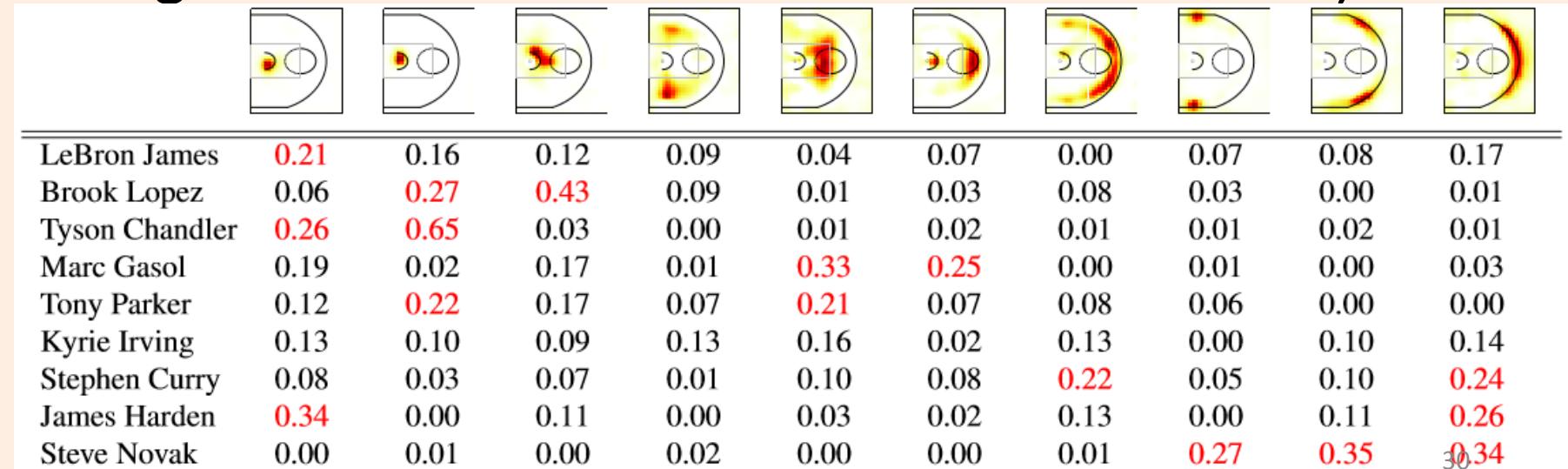
Then set negative values to 0.

# Application: Sports Analytics

- NBA shot charts:



- NMF (using “KL divergence” loss with k=10 and smoothed data).
  - Negative values would not make sense here.



# Application: Topic Modeling

- You have ‘n’ documents, ‘d’ bag-of-word features, want to find “topics”
- You can use NMF for this!
  - Interpretation of W: k topics, each with a selection of words
  - Interpretation of Z: each movie is a mixture of the k topics
- NMF makes much more sense than PCA
  - Each topic involves a small number of words
  - Each document has a small number of topics
- PCA would not make sense
  - you could have negative inclusion of a topic for a document
  - Topics can have negative words
  - all documents are a mixture of every possible topic and all topics involve every possible word
- So here we like both the sparsity and the non-negativity

# Regularized Matrix Factorization

- More recently people have considered L2-regularized PCA:

$$f(W, Z) = \frac{1}{2} \|ZW - X\|_F^2 + \frac{\lambda_1}{2} \|W\|_F^2 + \frac{\lambda_2}{2} \|Z\|_F^2$$

- Replaces normalization/orthogonality/sequential-fitting.
  - But requires regularization parameters  $\lambda_1$  and  $\lambda_2$ .
- Need to regularize W and Z because of scaling problem:
  - If you only regularize 'W' it doesn't do anything:
    - I could take unregularized solution, replace W by  $\alpha W$  for a tiny  $\alpha$  to shrink  $\|W\|_F$  as much as I want, then multiply Z by  $(1/\alpha)$  to get same solution.
  - Similarly, if you only regularize 'Z' it doesn't do anything.

# Sparse Matrix Factorization

- Instead of non-negativity, we could use L1-regularization:

$$f(W, Z) = \frac{1}{2} \|ZW - X\|_F^2 + \frac{\lambda_1}{2} \sum_{i=1}^n \|z_i\|_1 + \frac{\lambda_2}{2} \sum_{j=1}^d \|w_j\|_1$$

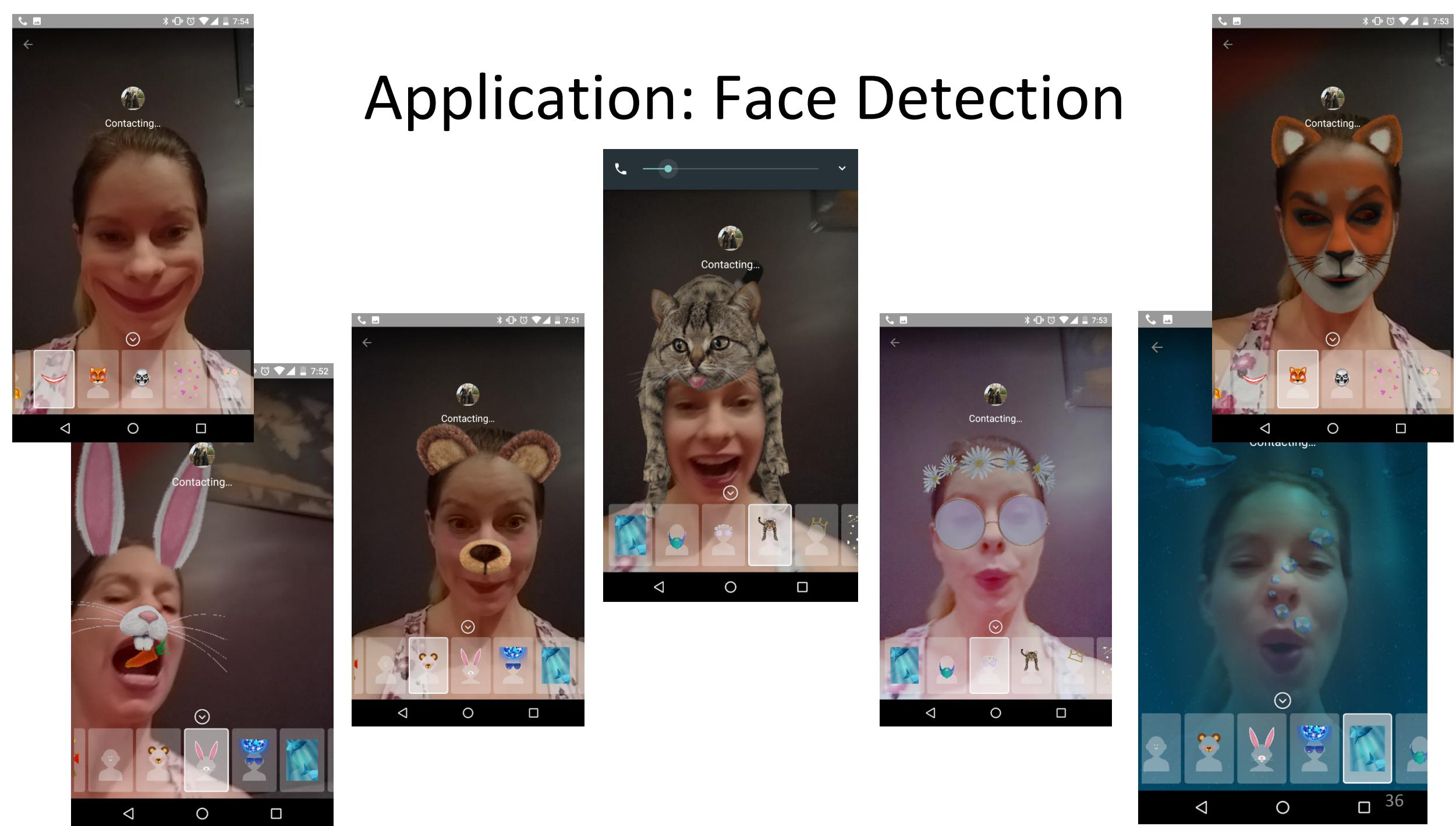
- Called **sparse coding** (L1 on ‘Z’) or **sparse dictionary learning** (L1 on ‘W’).
  - sklearn’s **SparsePCA** is L1 and ‘W’ and L2 on ‘Z’
- Disadvantage of using L1-regularization over non-negativity:
  - Sparsity controlled by  $\lambda_1$  and  $\lambda_2$  (so you need to set these)
- Advantage of using L1-regularization:
  - Sparsity controlled by  $\lambda_1$  and  $\lambda_2$  (so you can **control amount of sparsity**)
  - Also, negative coefficients often make sense.

# Sparsity: what is it good for?

- Sparsity: a vector/matrix with a bunch of zeros (often our ‘w’)
- Can be achieved in several ways:
  - Explicit feature selection, L1 regularization, non-negativity constraints
- Intuition: we want something “explained by a few factors”
  - NMF leads to **sparse Z and W**, whereas PCA does not.
- There can be big computational gains
  - We said earlier than SVM+kernels are fast because of the small number of support vectors. This has to do with **“sparsity in the dual”** (see CPSC 406)
- There are biological motivations
  - We believe there is “sparse coding” in the brain (few neurons in a pattern)
  - This might also mean more energy efficiency (both in the brain and in our tech)

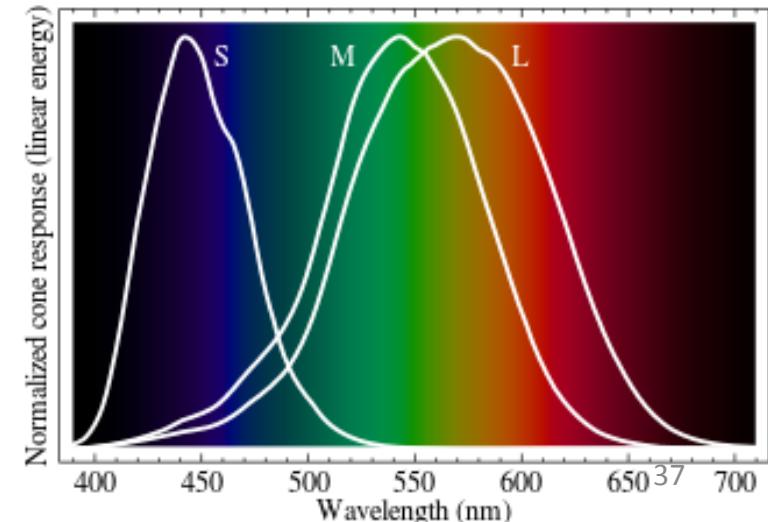
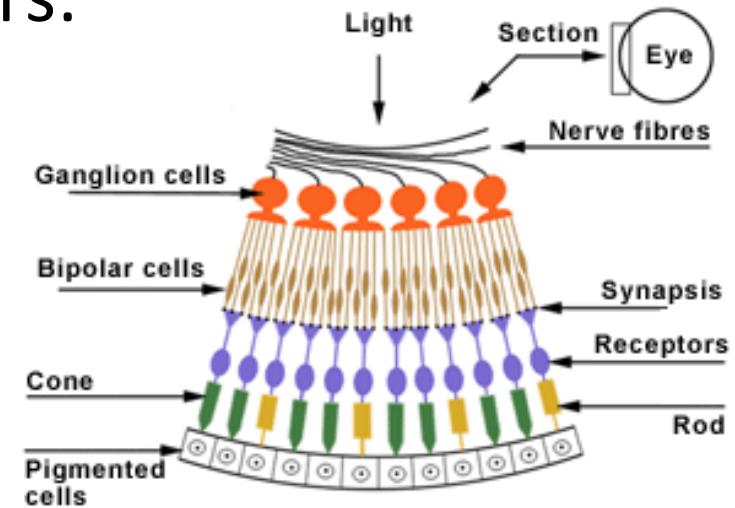
# Summary

- Non-negative matrix factorization leads to sparse ‘W’ and ‘Z’.
- Non-negativity constraints lead to sparse solution.
  - Projected gradient adds constraints to gradient descent.
- L1-regularization leads to other sparse latent-factor models.



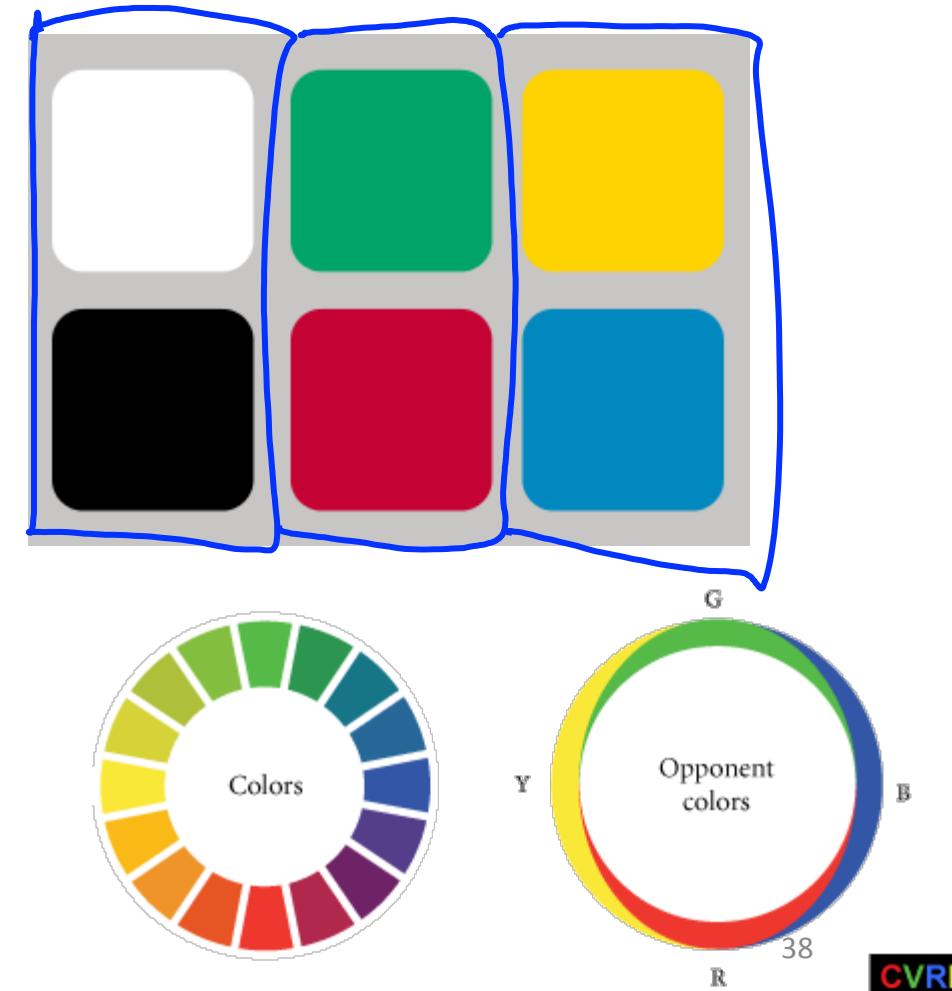
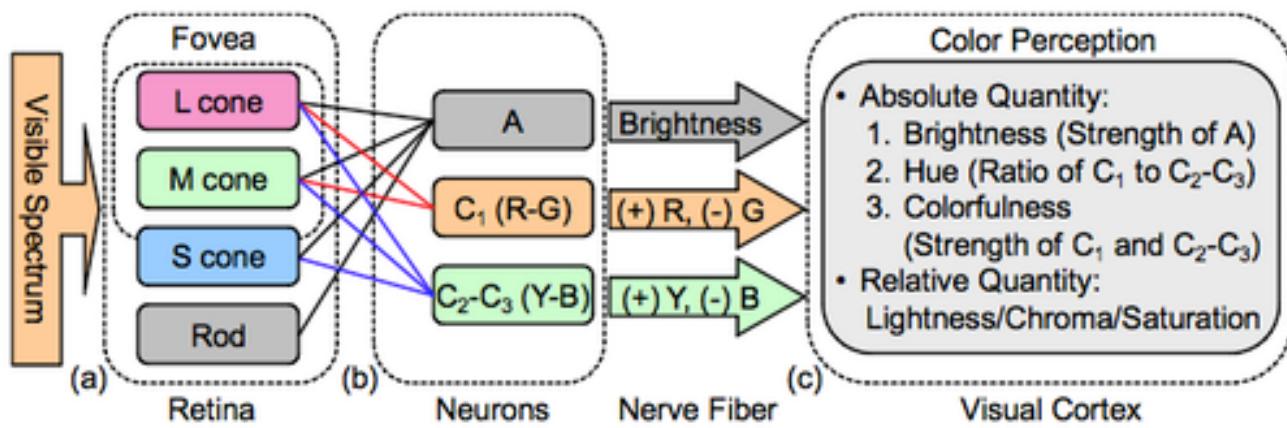
# Colour Opponency in the Human Eye

- Classic model of the eye is with 4 photoreceptors:
  - Rods (more sensitive to brightness).
  - L-Cones (most sensitive to red).
  - M-Cones (most sensitive to green).
  - S-Cones (most sensitive to blue).
- Two problems with this system:
  - Not orthogonal.
    - High correlation in particular between red/green.
  - We have 4 receptors for 3 colours.



# Colour Opponency in the Human Eye

- Bipolar and ganglion cells seem to code using “opponent colors”:– 3-variable orthogonal basis:



# Colour Opponency Representation

For this pixel,  
eye gets 4 signals



$$= w_1$$

First row  
of  $W$

(First PC)

Analogous to means in k-means.



brightness

Can represent 4 original values with  
these 3  $z_i$  values and  
matrix ' $W$ '



red/green

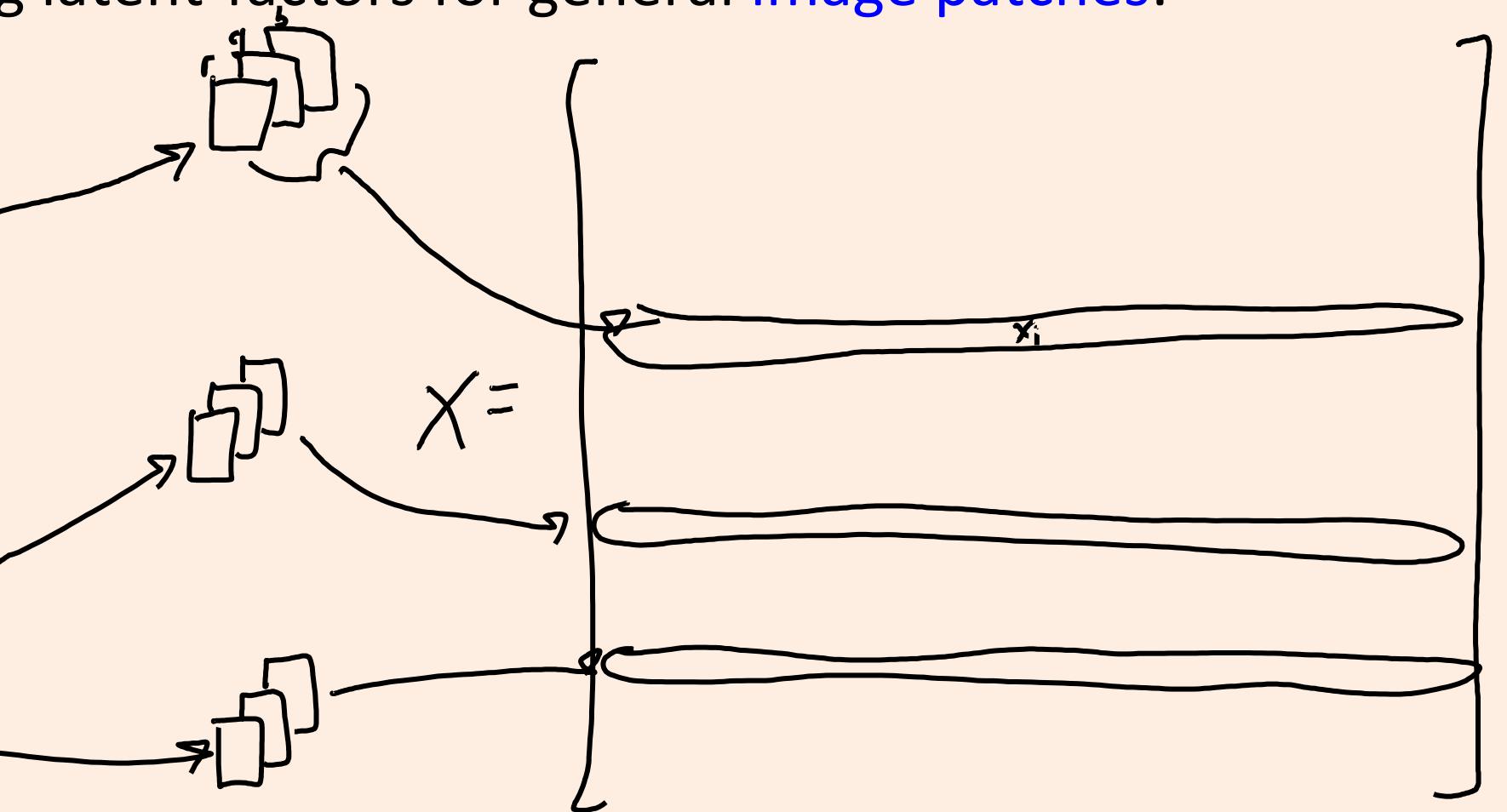
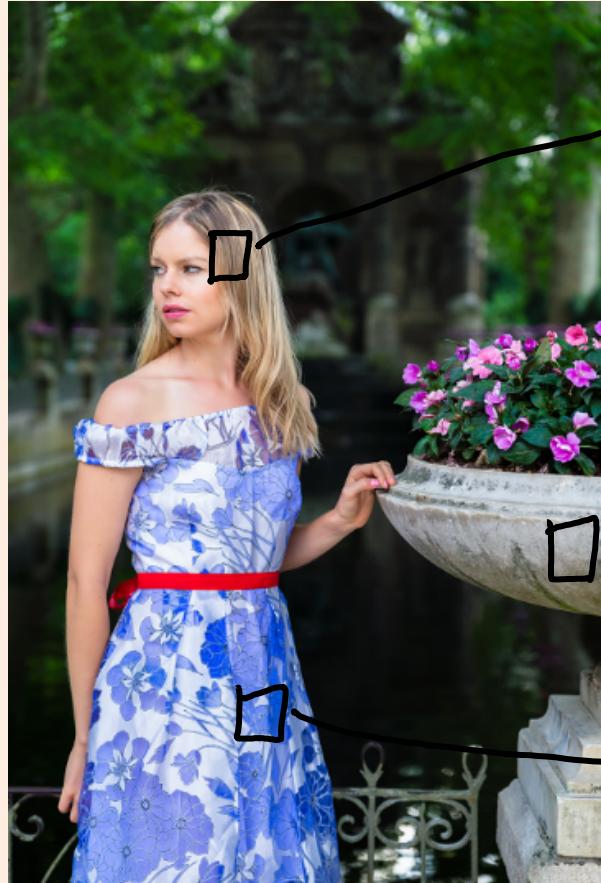
+  $w_2$   
Second  
row  
 $(4 \times 1)$



blue/yellow

# Latent-Factor Models for Image Patches

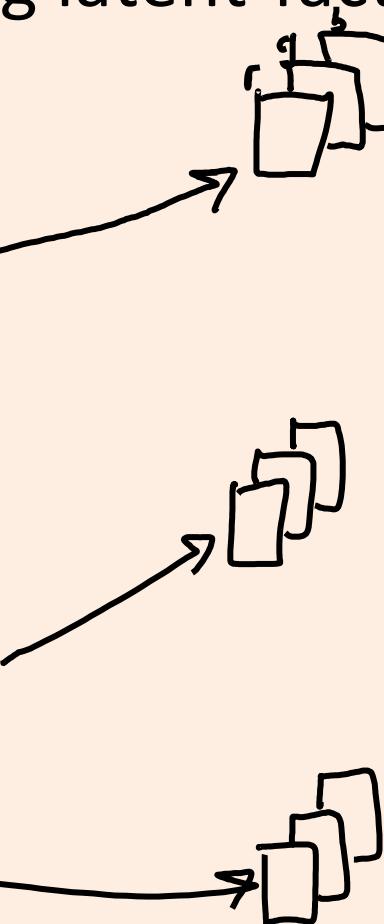
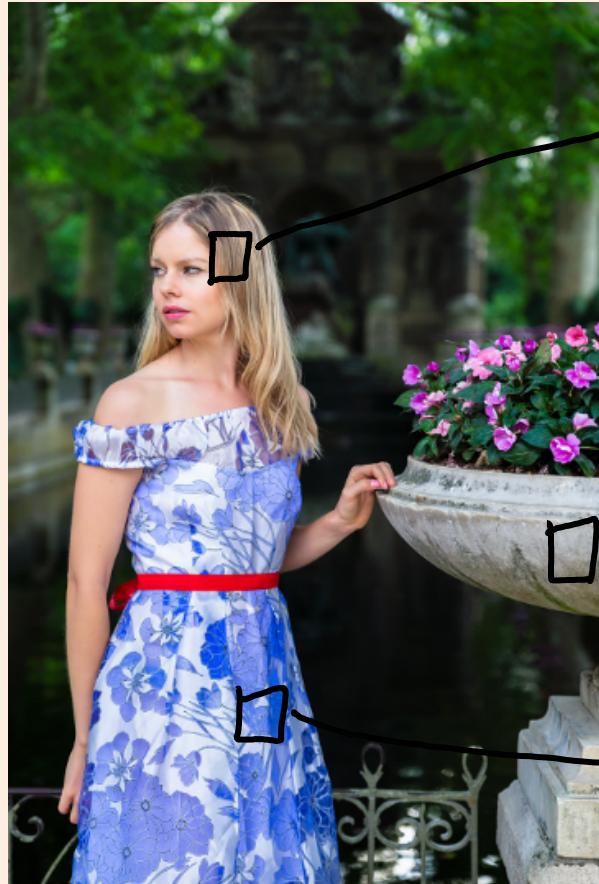
- Consider building latent-factors for general **image patches**:



Size of  $X$ : (image height)  $\times$  (image width) by ((patch height)  $\times$  (patch width)  $\times$  3)

# Latent-Factor Models for Image Patches

- Consider building latent-factors for general **image patches**:



Typical pre-processing:

1. Usual variable centering
2. “Whiten” patches.  
(remove correlations)

# Application: Image Restoration

The Salinas Valley is in Northern California. It is a long narrow valley between two ranges of mountains, and the Salinas River winds and twists as it center until it falls at last into Monterey Bay.

remember my childhood names for grasses and desert flowers. I remember where a road may live and what, like the birds, passed in the summer and what trees and seasons smelled like how people looked and walked and sounded even. The memory of stories is very rich.

I remember that the Gabilan Range to the east of the valley were big gray mountains full of bus and boulders and a kind of mistiness, so that you wanted to climb into their warm bottoms almost as you want to climb into the top of a beloved mother. They were backbone mountains with a brown grass love. The same bushes stood up against the sky to the west and kept the valley from the open sea, and they were dark and growing-unfriendly and dangerous. I always found in myself a great as well as a love of east. Where I ever got such an idea I cannot say, unless it could be that the morning came over the peaks of the Gabilans and the night dimmed back from the ridges of the Santa Lucia. It may be that the birth and death of the day had some part in my feeling about the two ranges of mountains.

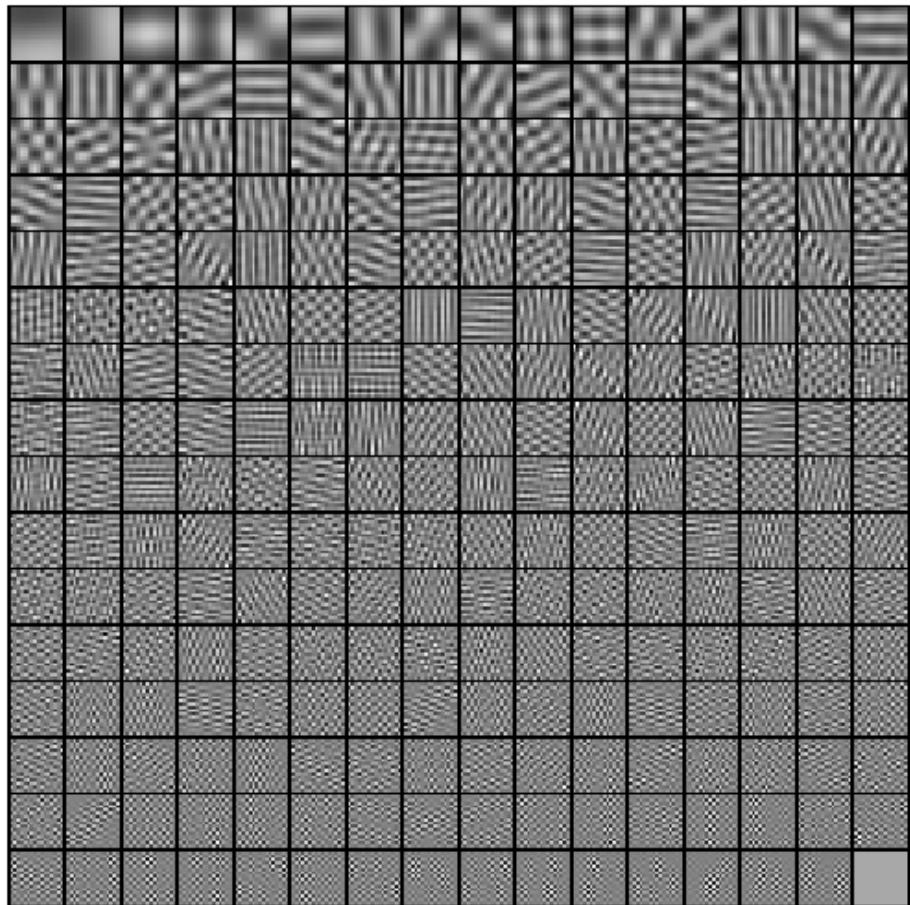
From both sides of the valley little streams slipped out of the hill canyons and falling the bed of the Salinas River. In the winter of wet years the streams overflowed and they swelled the river until sometimes it ranged and boiled, bank full, and though it was a destroyer. The river, tore the edges of the farm lands and washed whole acres down, it loosed barns and houses into itself, to go Roaring and bubbling away. It trapped cows and pigs and sheep and drowned them in its muddy green water, and carried them to the sea. Then when the tide turned back, the river dried up, and the mud was left, and the flooded fields appeared, and so the animals, the sheep, dogs, and all other forms, went about what was left in the deep water places under a king's mantle. The fishes and snakes, also, had to walk, or swim, or burrow, or roll, in their upper blossoms. The Salinas was only a part-time river. The winter sun drove it underground, it was not a river then at all, but it was the only river we had and not to be boasted about. It was dangerous! It was in a wet winter and had dry flowers in dry summer. You can boast about anything if it's all you have. Maybe the less you have, the more you are required to boast.

The floor of the Salinas Valley, between the ridges and below the foothills, is level because this valley used to be the bottom of a hundred-mile inlet from the sea. The lower mouth of Moss Landing Bay exists now ago the entrance to this long inland water since, fifty miles down the valley, my father built a wall. The dam stands up first with logs and then with gravel and then with white sea sand fused shells and even pebbles.

ing mountains  
alley from th  
If a dread o  
morning cam



# Latent-Factor Models for Image Patches

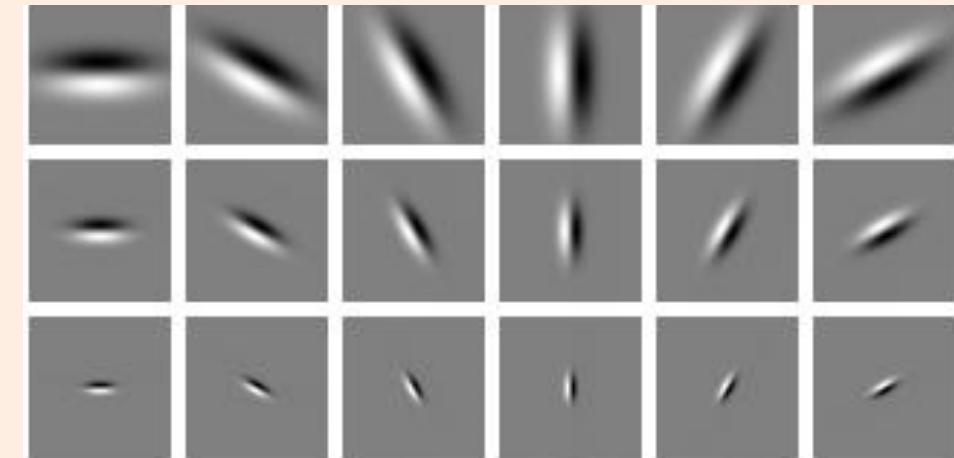


(b) Principal components.

Orthogonal bases don't seem right:

- Few PCs do almost everything.
- Most PCs do almost nothing.

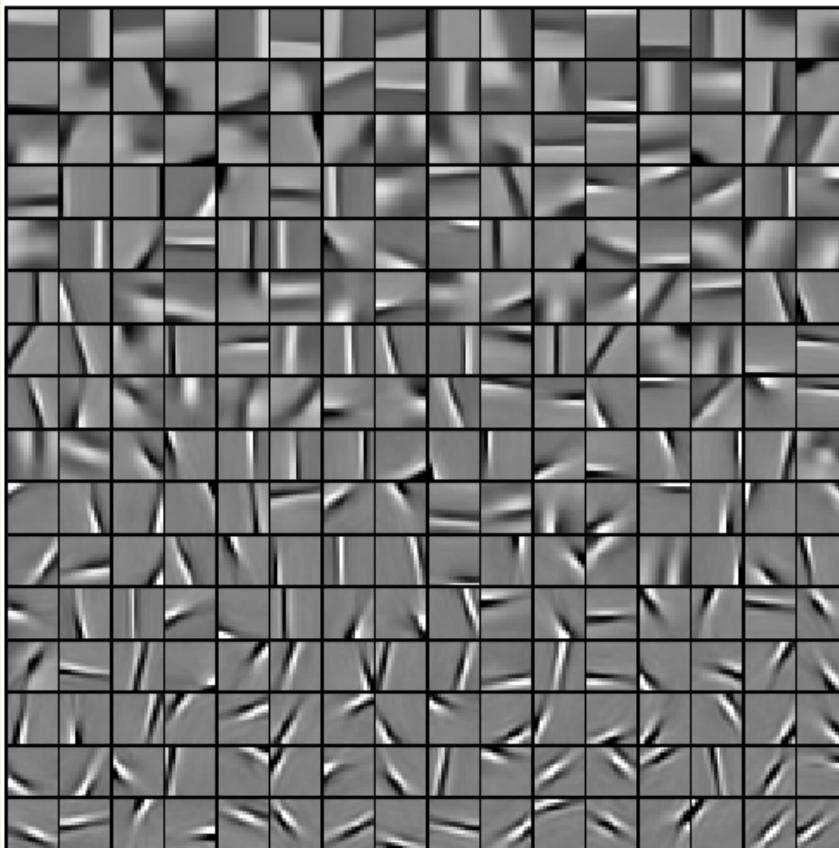
We believe “simple cells” in visual cortex use:



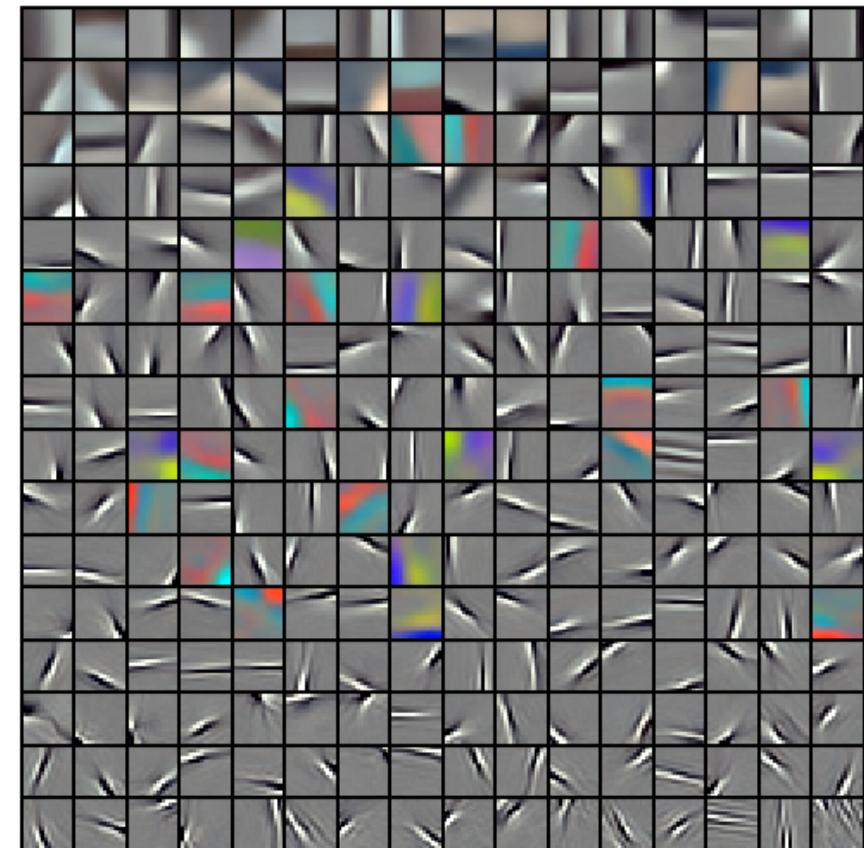
‘Gabor’ filters

# Latent-Factor Models for Image Patches

- Results from a sparse (non-orthogonal) latent factor model:



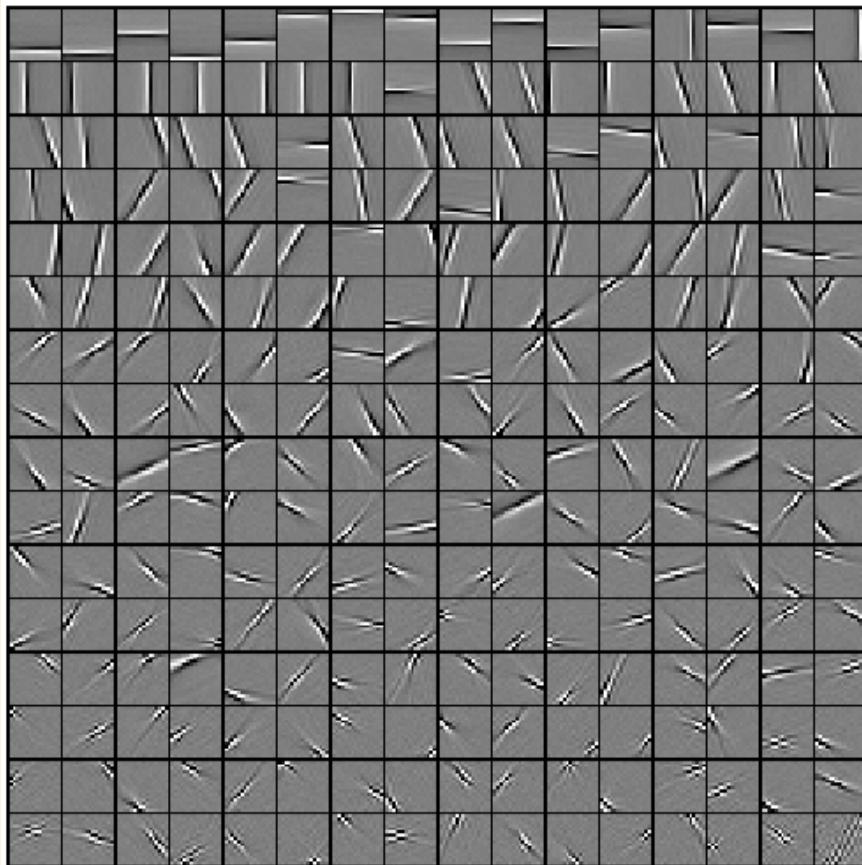
(a) With centering - gray.



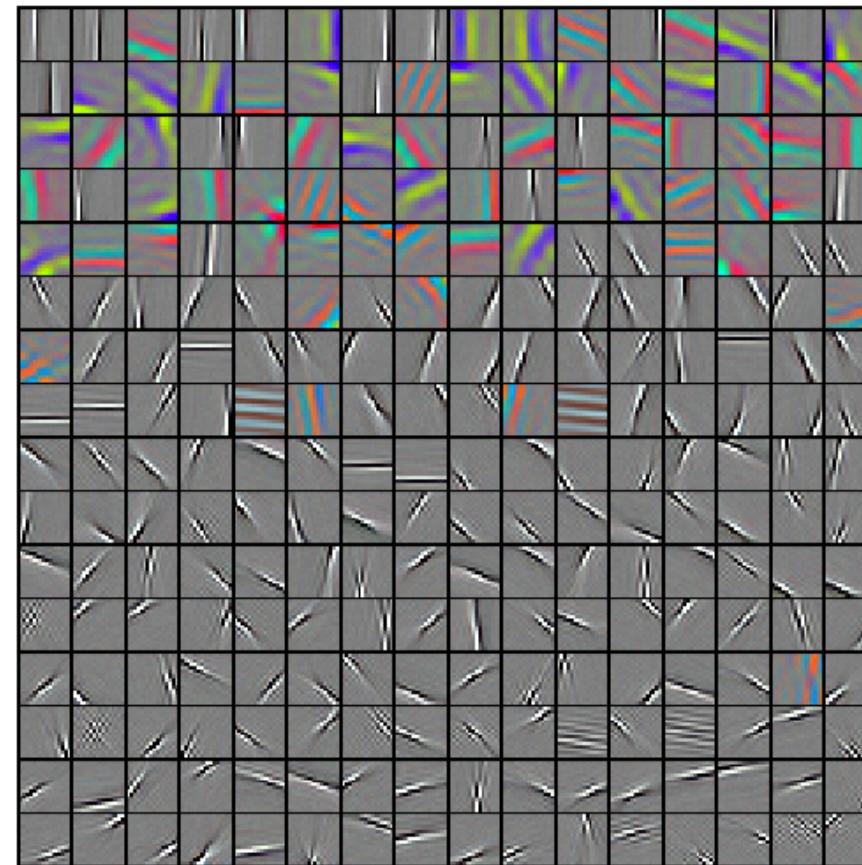
(b) With centering - RGB.

# Latent-Factor Models for Image Patches

- Results from a “sparse” (non-orthogonal) latent-factor model:



(c) With whitening - gray.

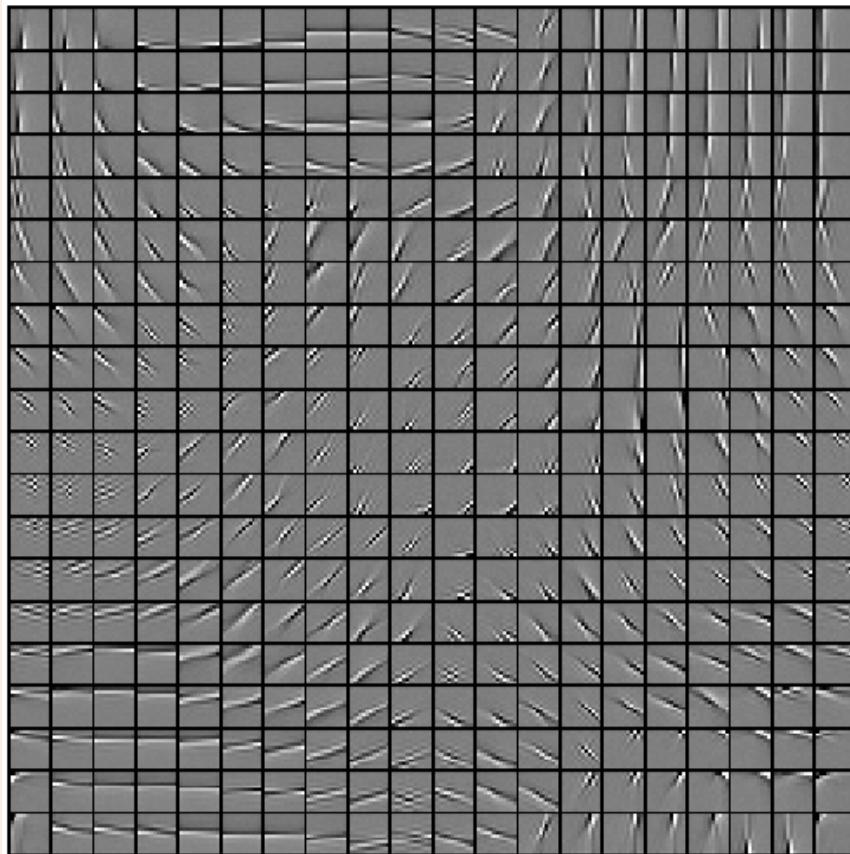


(d) With whitening - RGB.

"colour  
opponency"

# Recent Work: Structured Sparsity

- Basis learned with a variant of “structured sparsity”:

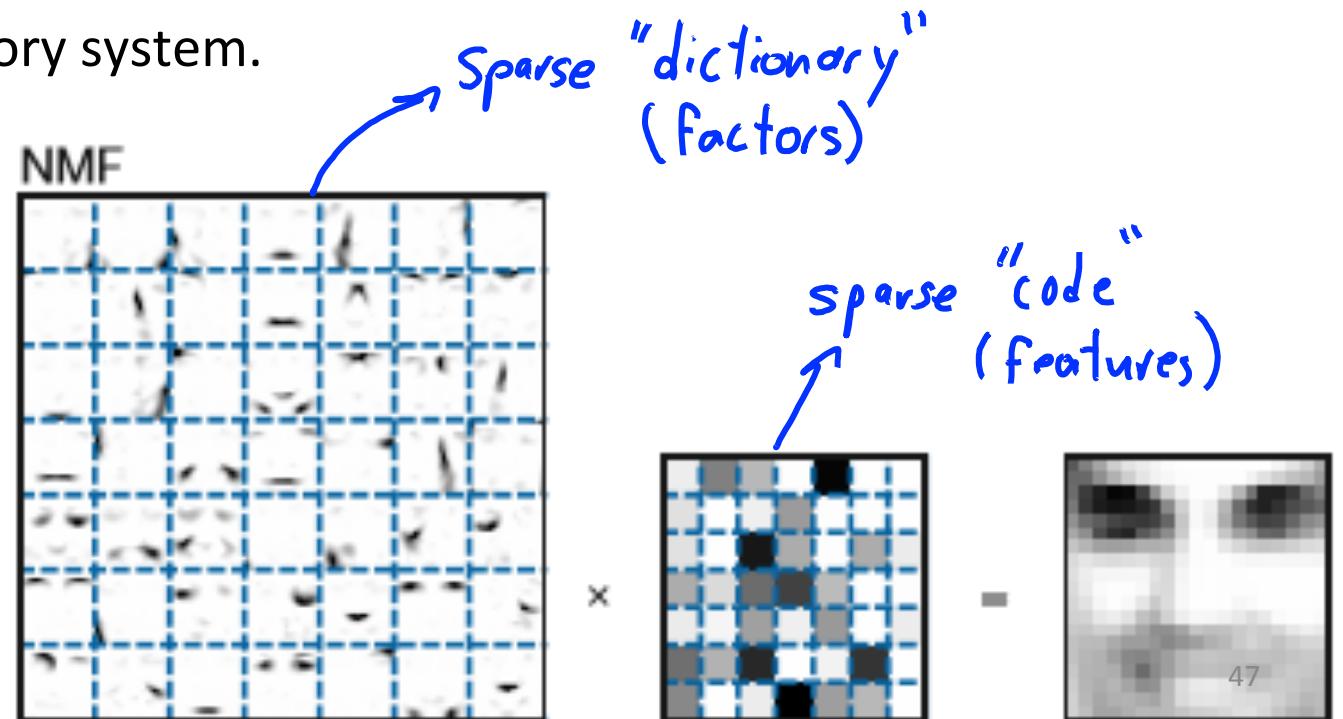


(b) With  $4 \times 4$  neighborhood.

Similar to “cortical columns”  
theory in visual cortex.

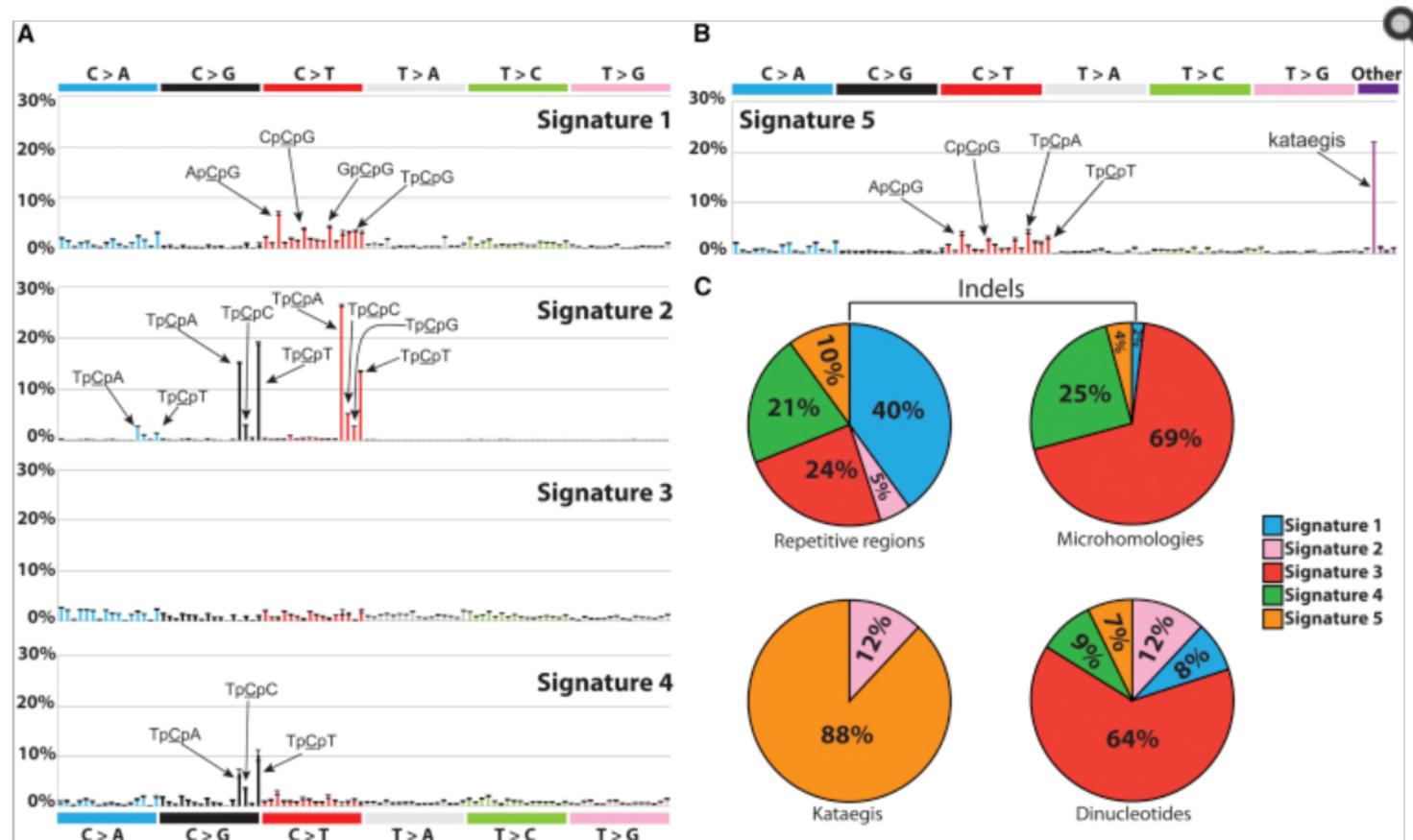
# Representing Faces

- Why sparse coding?
  - “Parts” are intuitive, and brains seem to use sparse representation.
  - Energy efficiency if using sparse code.
  - Increase number of concepts you can memorize?
    - Some evidence in fruit fly olfactory system.



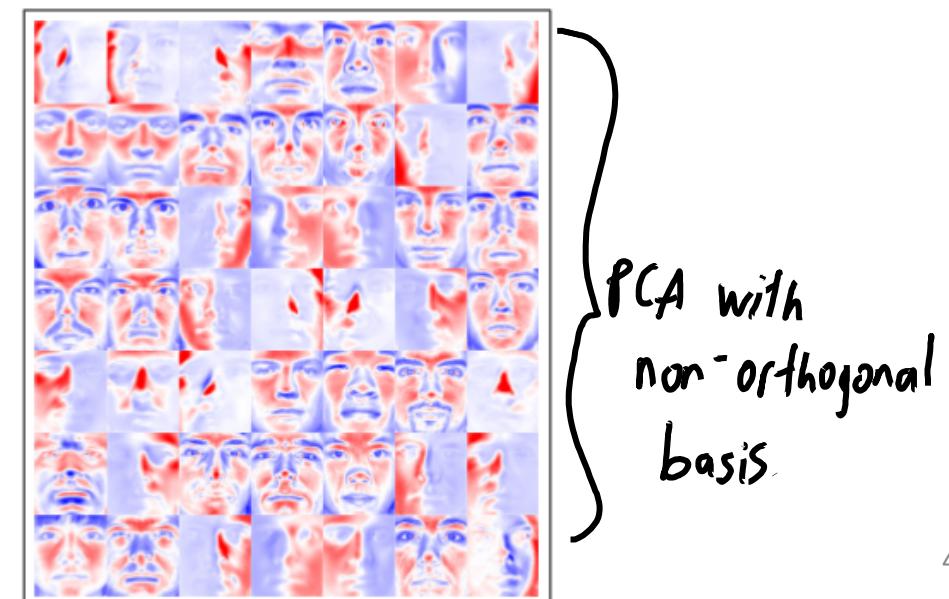
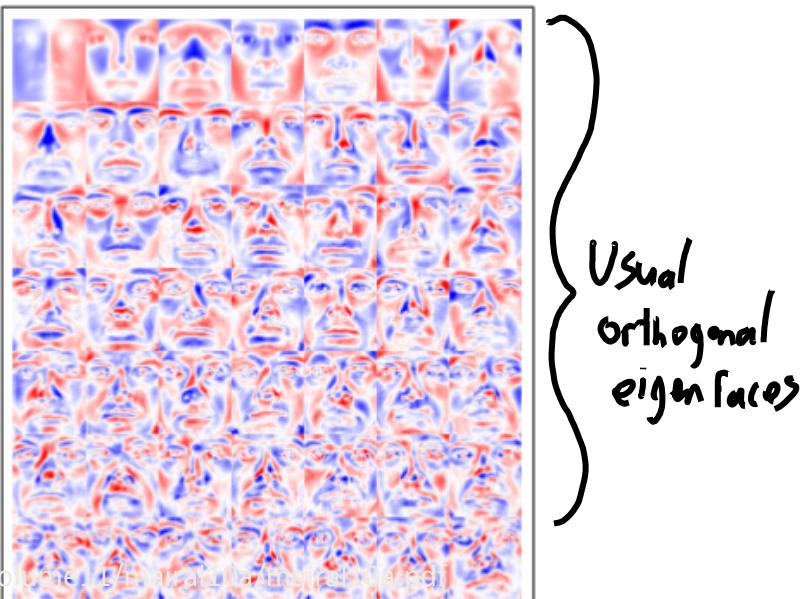
# Application: Cancer “Signatures”

- What are common sets of mutations in different cancers?
  - May lead to new treatment options.



# Regularized Matrix Factorization

- For many PCA applications, ordering orthogonal PCs makes sense.
  - Latent factors are independent of each other.
  - We definitely want this for visualization.
- In other cases, ordering orthogonal PCs doesn't make sense.
  - We might not expect a natural “ordering”.



# Sparse Matrix Factorization

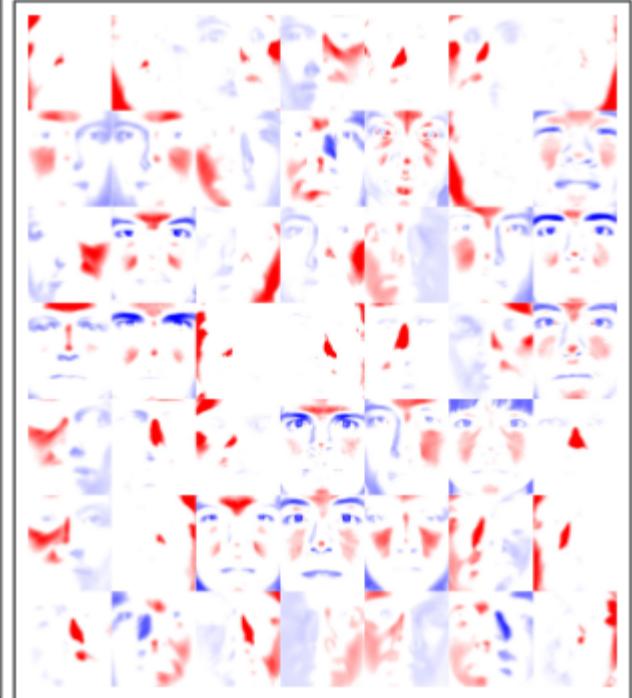
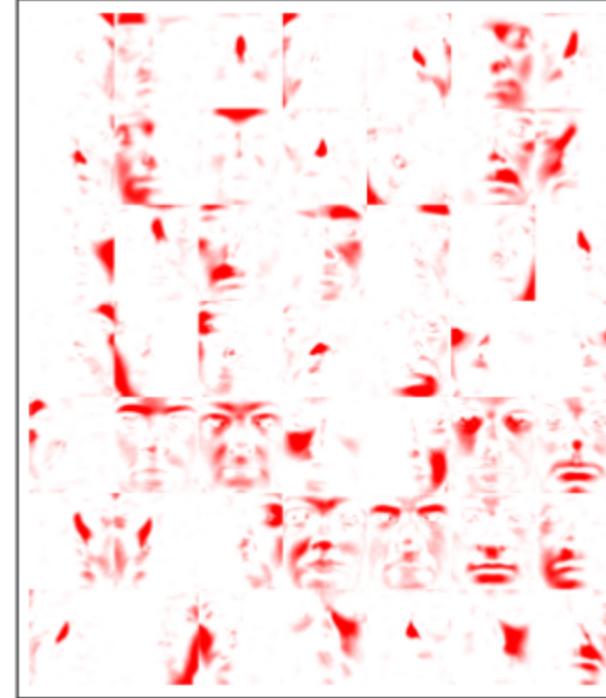
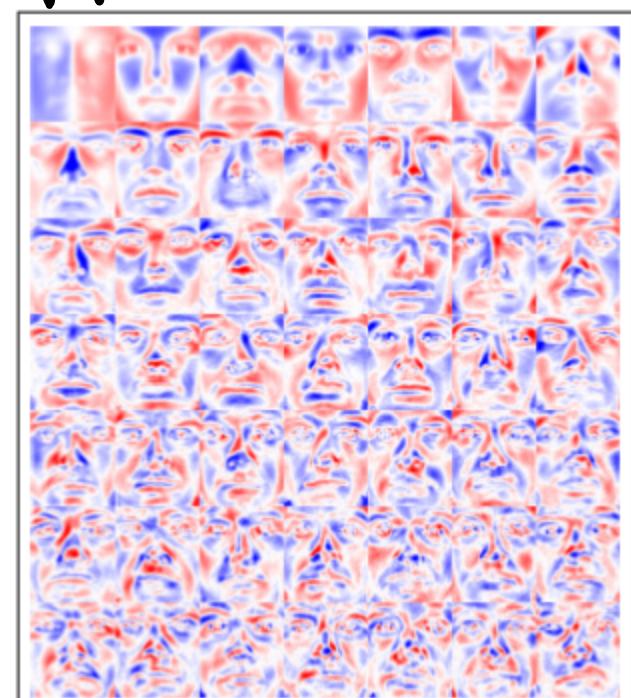
- Instead of non-negativity, we could use L1-regularization:

$$f(W, Z) = \frac{1}{2} \|ZW - X\|_F^2 + \frac{\lambda_1}{2} \sum_{i=1}^n \|z_i\|_1 + \frac{\lambda_2}{2} \sum_{j=1}^d \|w_j\|_1$$

- Called **sparse coding** (L1 on ‘Z’) or **sparse dictionary learning** (L1 on ‘W’).
- Many variations exist:
  - Mixing L2-regularization and L1-regularization.
    - Or normalizing ‘W’ (in L2-norm or L1-norm) and regularizing ‘Z’.
  - **K-SVD** constrains each  $z_i$  to have at most ‘k’ non-zeroes:
    - K-means is special case where  $k = 1$ .
    - PCA is special case where  $k = d$ .

# Matrix Factorization with L1-Regularization

blue: negative  
red: positive



PCA without orthogonality

sparsity due to  
non-negativity

sparsity due to  
L<sub>1</sub>-regularization

# Recent Work: Structured Sparsity

- “Structured sparsity” considers dependencies in sparsity patterns.
  - Can enforce that “parts” are convex regions.

