# CPSC 340:
# Machine Learning and Data Mining

Sparse Matrix Factorization
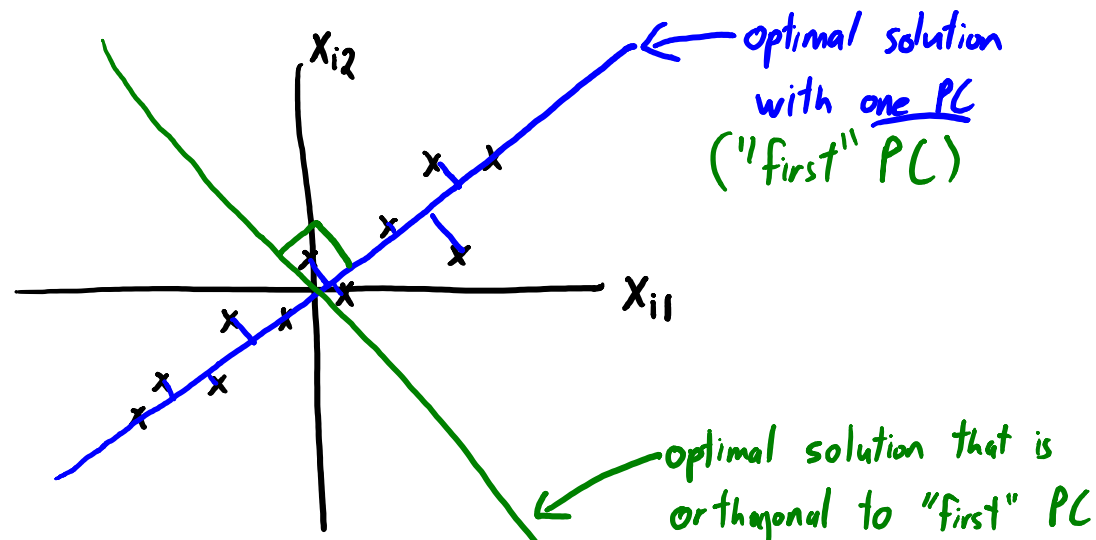
Fall 2017

# Admin

- **Assignment 4**:
  - Due Friday.

- **Assignment 5**:
  - Posted, due Monday of last week of classes

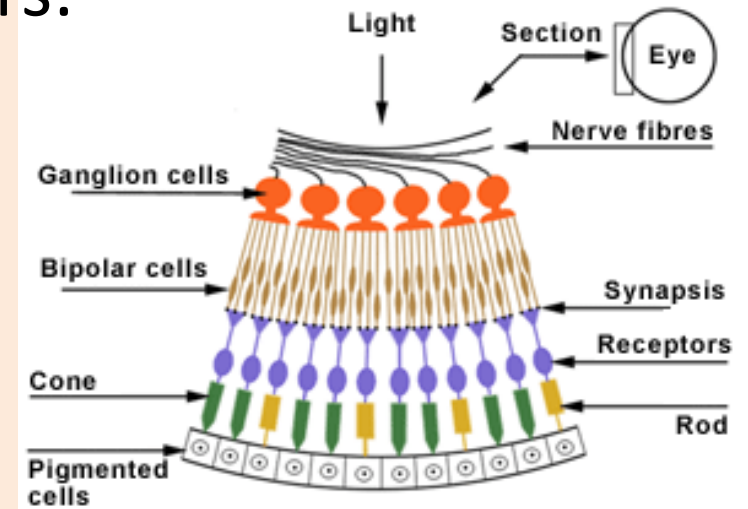# Last Time: PCA with Orthogonal/Sequential Basis

- 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'.

$$\|w_c\| = 1 \quad \text{and} \quad w_c^{\top} w_{c'} = 0 \quad \text{for} \quad c' \neq c$$

  - And fit the rows in order:

    - First row "explains the most variance" or "reduces error the most".
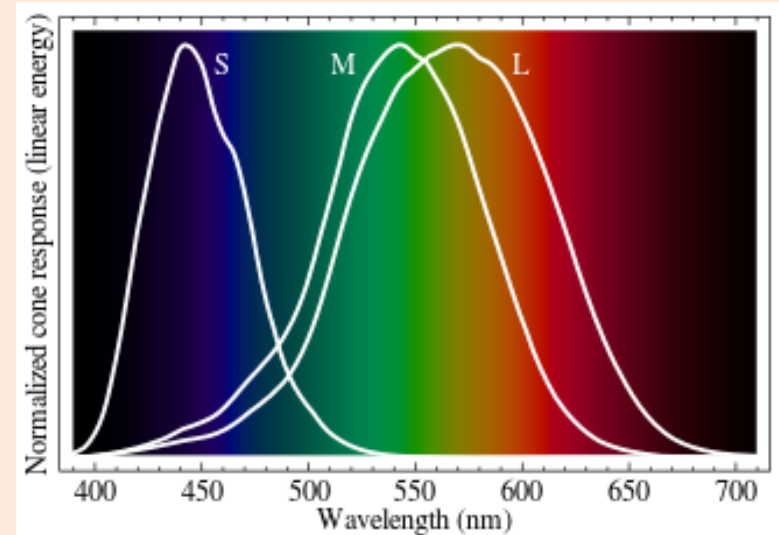
# Colour Opponency in the Human Eye

- 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).
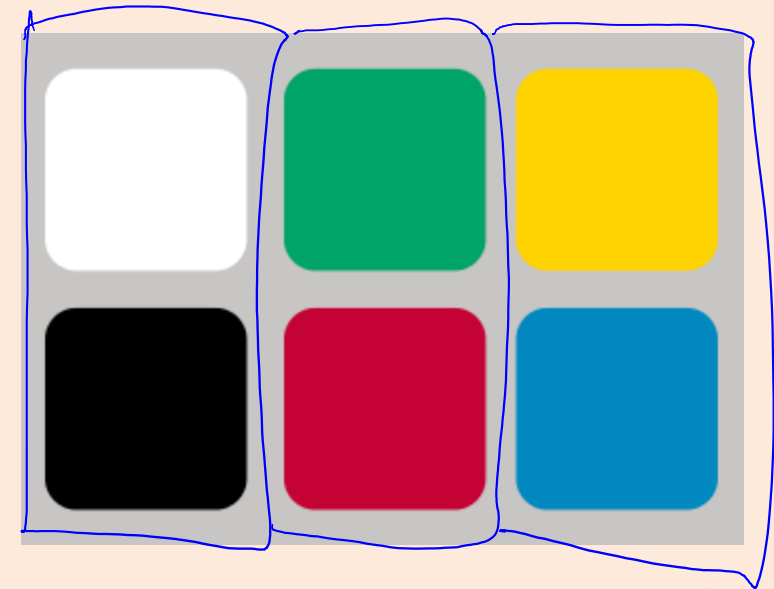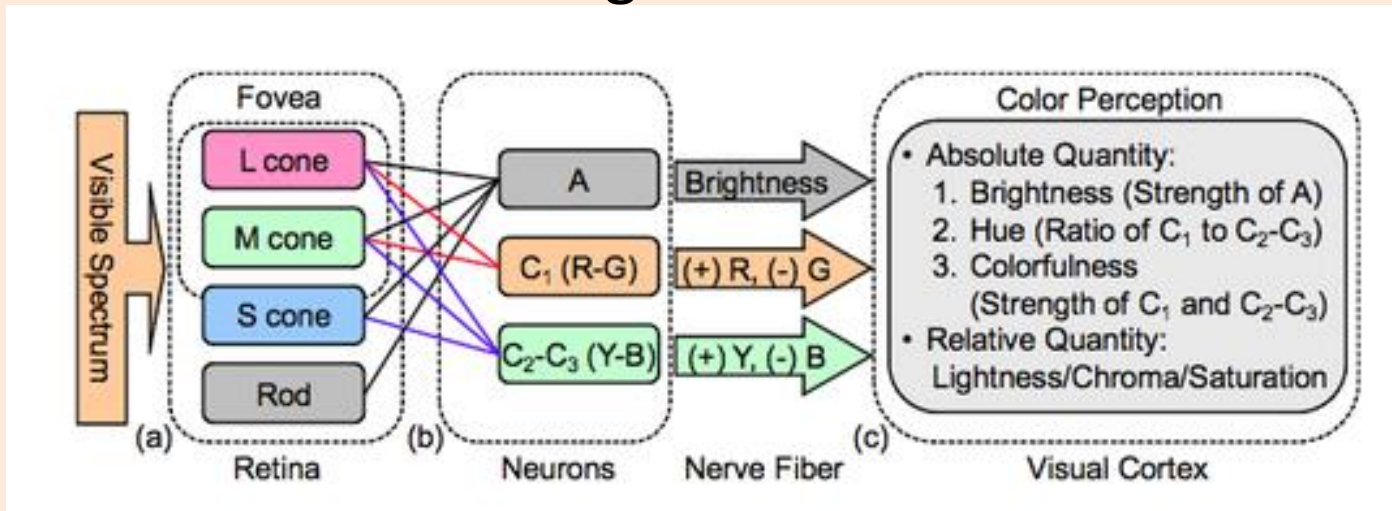  - S-Cones (most sensitive to blue).

- Two problems with this system:
  - Not orthogonal.
    - High correlation in particular between red/green.
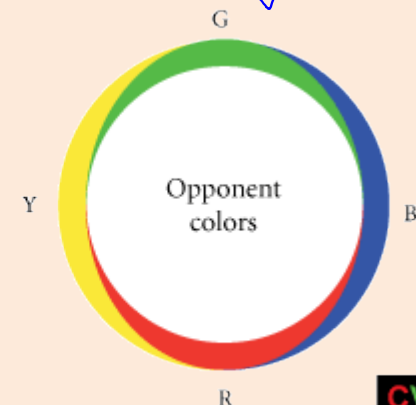  - We have 4 receptors for 3 colours.

# Colour Opponency in the Human Eye

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



- This is similar to PCA (d = 4, k = 3).

# Colour Opponency Representation



For this pixel, eye gets 4 signals

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

$= W_1$

First row of W (First PC)

Analogous to means in k-means.

brightness

$+ W_2$

Second row (4×1)

red/green
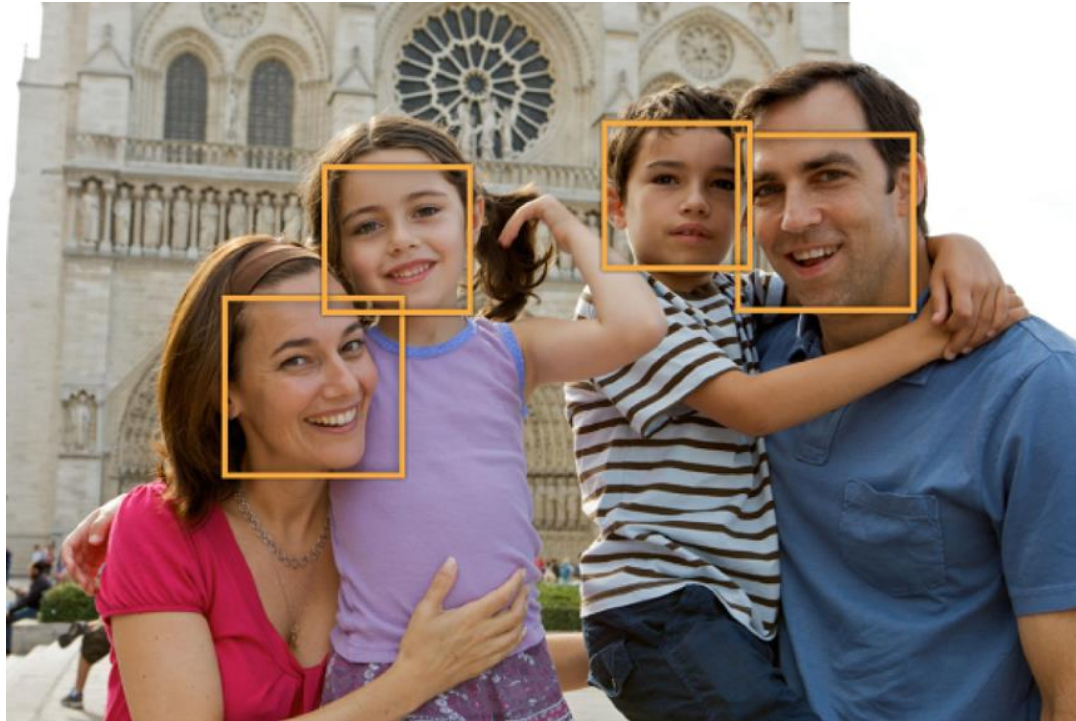
$+ W_3$
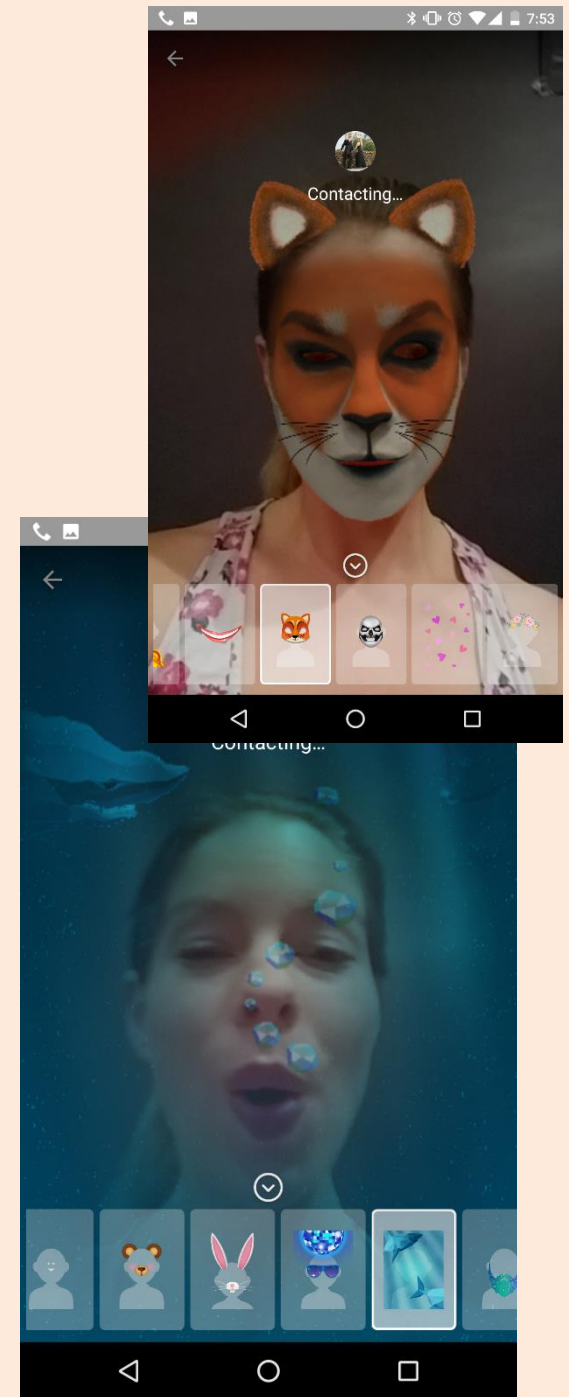
Third row (4×1)

blue/yellow

# Application: Face Detection

- Consider problem of face detection:



- Classic methods use "eigenfaces" as basis:
  - PCA applied to images of faces.

# Application: Face Detection

# 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.



Replace each $x_{ij}$ by $x_{ij} - \mu_j$
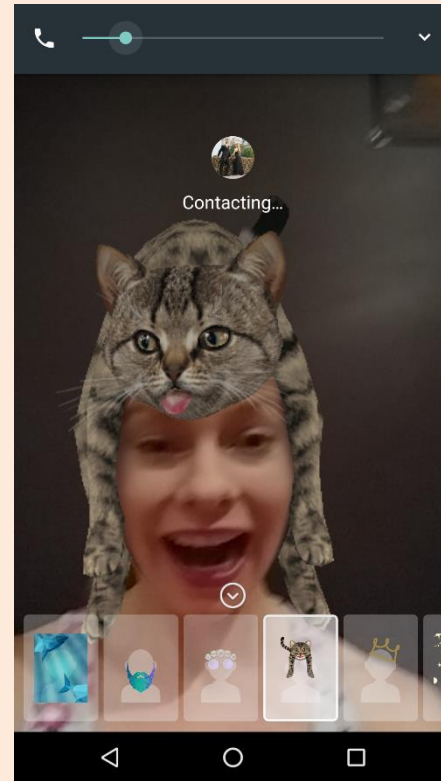
Each row of $X$ will be pixels in __one__ image:

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

# Eigenfaces

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



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

# Eigenfaces

Compute top 'k' PCs on centered data:



PC2

PC1

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:

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

$\hat{x}_i$    $\mu$    PC1 (first row of W)    PC2    PC3

# Eigenfaces

106 of the original faces:



"Eigenface" representation:

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

PC1 (first row of $W$)

# Eigenfaces

Reconstruction with $k=0$



Variance explained: $0\%$

"Eigenface" representation:



$\hat{x_i}$ $=$ $\mu$ $+z_{i1}$ PC1 $+z_{i2}$ PC2 $+z_{i3}$ PC3 $+\cdots$

(first row of W)

# Eigenfaces

Reconstruction with $k=1$



Variance explained: 34%

PCA Visualization:



$$z_i = w_c^T x_i$$

"Eigenface" representation:



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

(first row of W)

# Eigenfaces

Reconstruction with $k=2$



Variance explained: $71\%$

PCA Visualization:



"Eigenface" representation:

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

(first row of $W$)

# Eigenfaces

Reconstruction with $k=3$



Variance explained: $76\%$

PCA Visualization:
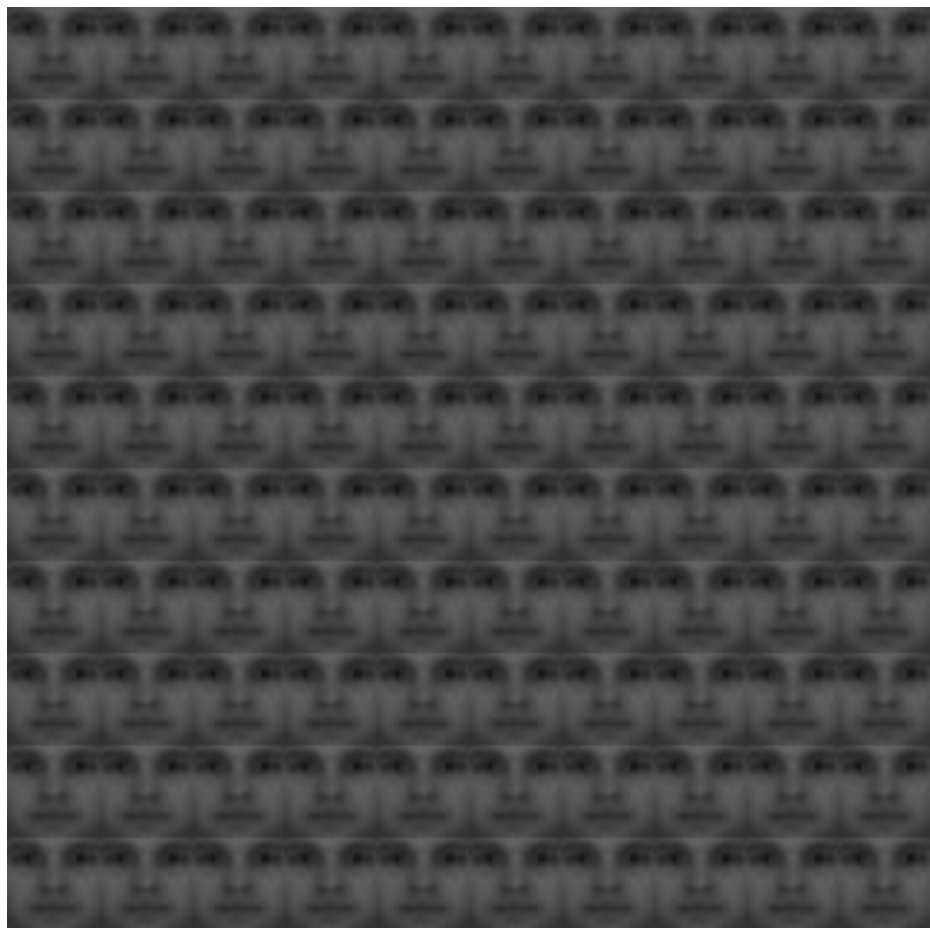


"Eigenface" representation:



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

(first row of $W$)

# 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:



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

Plus these "eigenfaces" and the mean.

# VQ vs. PCA vs. NMF

- But how *should* we represent faces?
  - Vector quantization (k-means).
    - Replace face by the average face in a cluster.
    - 'Grandmother cell': one neuron = one face.
    - Can't distinguish between people in the same cluster (only 'k' possible faces).
    - Almost certainly not true: too few neurons.



$$\hat{X}_i = z_{i1} * w_1 + z_{i2} * w_2 + z_{i3} * w_3 + z_{i4} * w_4 + z_{i5} * w_5 + z_{i6} * w_6$$

# VQ vs. PCA vs. NMF

- But *should* we represent faces?
  - Vector quantization (k-means).
  - PCA (orthogonal basis).
    - Global average plus linear combination of "eigenfaces".
    - "Distributed representation".
      - Coded by pattern of group of neurons: can represent infinite number of faces by changing $z_i$.
    - But "eigenfaces" are not intuitive ingredients for faces.
      - PCA tends to use positive/negative cancelling bases.



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

# 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.



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

# 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.



Sparse "dictionary" (factors)

sparse "code" (features)

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

- Consider our usual least squares problem:

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

- 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:
  - No physical interpretation to negative weights.
  - If $x_{ij}$ is amount of product you produce, what does $w_j < 0$ mean?

- Non-negativity leads to sparsity...

# Sparsity and Non-Negative Least Squares

- Consider 1D non-negative least squares objective:

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

- 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:

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

- 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.

$$\text{Compute } w = (X^\top X) \backslash (X^\top y)$$

$$\text{Set } \quad w_j = \max\{0, w_j\}$$

  - This is correct when d = 1.
  - Can be worse than setting w = 0 when d ≥ 2.

# 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:

    $$w^{t+\frac{1}{2}} = w^t - \alpha^t \nabla f(w^t)$$

    - After each step, set negative values to 0.

    $$w_j^{t+1} = \max\left\{0, w_j^{t+\frac{1}{2}}\right\}$$

    - 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:

$$w^{t+\frac{1}{2}} = w^t - \alpha^t \nabla f(w^t) \qquad w_j^{t+1} = \max\{0, w_j^{t+\frac{1}{2}}\}$$

  - 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.m)

# Projected-Gradient for NMF

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

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

  – 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.

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

(keep other values of W and Z fixed)

  – Non-convex and (unlike PCA) is sensitive to initialization.
    - Hard to find the global optimum.
    - Typically use random initialization.

Then set negative values to 0.

# Application: Sports Analytics

- NBA shot charts:



Stephen Curry (940 shots)   LeBron James (315 shots)

- NMF (using "KL divergence" loss with k=10 and smoothed data).

  - Negative values would not make sense here.



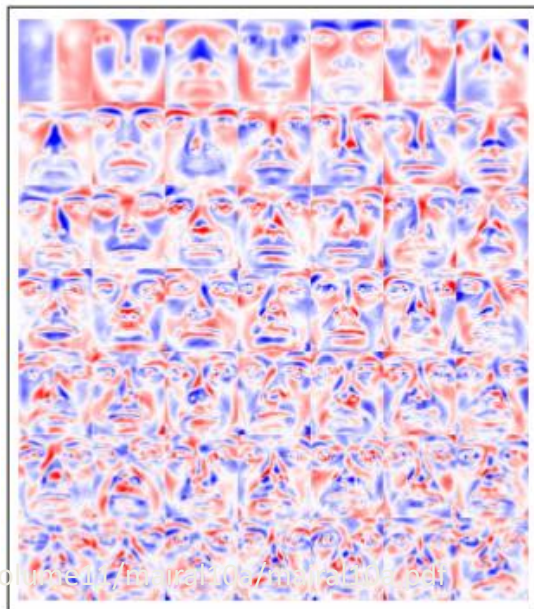| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| LeBron James | 0.21 | 0.16 | 0.12 | 0.09 | 0.04 | 0.07 | 0.00 | 0.07 | 0.08 | 0.17 |
| Brook Lopez | 0.06 | 0.27 | 0.43 | 0.09 | 0.01 | 0.03 | 0.08 | 0.03 | 0.00 | 0.01 |
| Tyson Chandler | 0.26 | 0.65 | 0.03 | 0.00 | 0.01 | 0.02 | 0.01 | 0.01 | 0.02 | 0.01 |
| Marc Gasol | 0.19 | 0.02 | 0.17 | 0.01 | 0.33 | 0.25 | 0.00 | 0.01 | 0.00 | 0.03 |
| Tony Parker | 0.12 | 0.22 | 0.17 | 0.07 | 0.21 | 0.07 | 0.08 | 0.06 | 0.00 | 0.00 |
| Kyrie Irving | 0.13 | 0.10 | 0.09 | 0.13 | 0.16 | 0.02 | 0.13 | 0.00 | 0.10 | 0.14 |
| Stephen Curry | 0.08 | 0.03 | 0.07 | 0.01 | 0.10 | 0.08 | 0.22 | 0.05 | 0.10 | 0.24 |
| James Harden | 0.34 | 0.00 | 0.11 | 0.00 | 0.03 | 0.02 | 0.13 | 0.00 | 0.11 | 0.26 |
| Steve Novak | 0.00 | 0.01 | 0.00 | 0.02 | 0.00 | 0.00 | 0.01 | 0.27 | 0.35 | 0.34 |

# Application: Cancer "Signatures"

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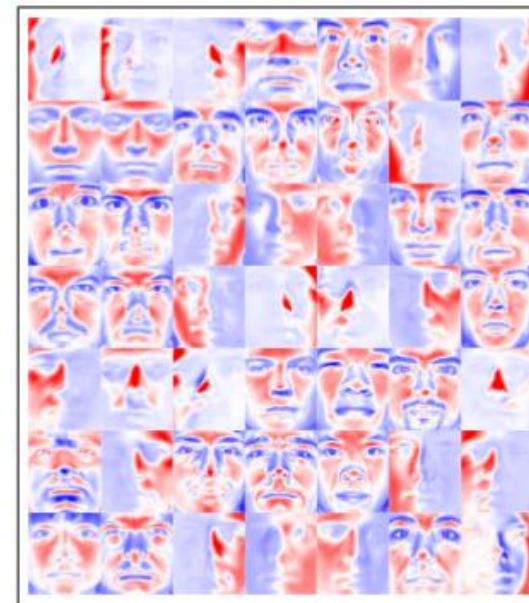  - May lead to new treatment options.

# 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".



Usual orthogonal eigen faces

PCA with non-orthogonal basis.

# Regularized Matrix Factorization

- More recently people have considered <span style="color:blue">L2-regularized PCA</span>:

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

- <span style="color:green">Replaces normalization/orthogonality/sequential-fitting</span>.
  - But requires <span style="color:red">regularization parameters</span> $\lambda_1$ and $\lambda_2$.

- <span style="color:green">Need to regularize W and Z</span> because of scaling problem:
  - <span style="color:red">If you only regularize 'W' it doesn't do anything</span>:
    - 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.
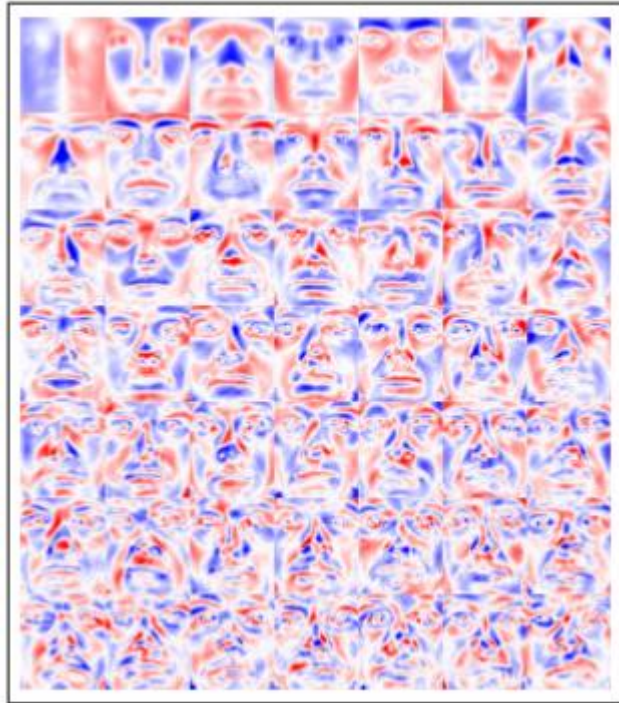
# Sparse Matrix Factorization

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

$$f(W, Z) = \frac{1}{2}\| ZW - X \|_F^2 + \frac{\lambda_1}{2}\sum_{i=1}^{n}\|z_i\|_1 + \frac{\lambda_2}{2}\sum_{j=1}^{d}\|w_j\|_1$$
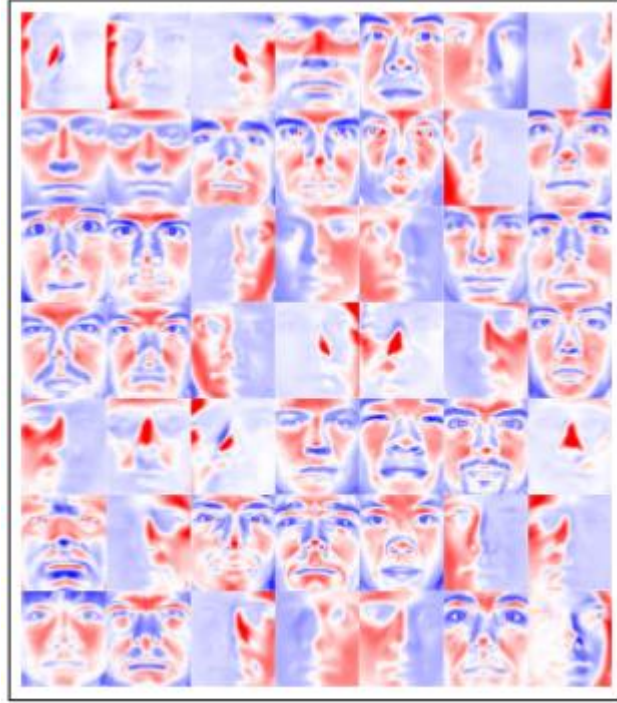
  – Called sparse coding (L1 on 'Z') or sparse dictionary learning (L1 on 'W').

- 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:
  – Negative coefficients usually make sense.
  – Sparsity controlled by $\lambda_1$ and $\lambda_2$, so you can control amount of sparsity.

# 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



(a) PCA

(e) Dictionary Learning

(c) NMF

(d) SPCA, $\tau = 30\%$

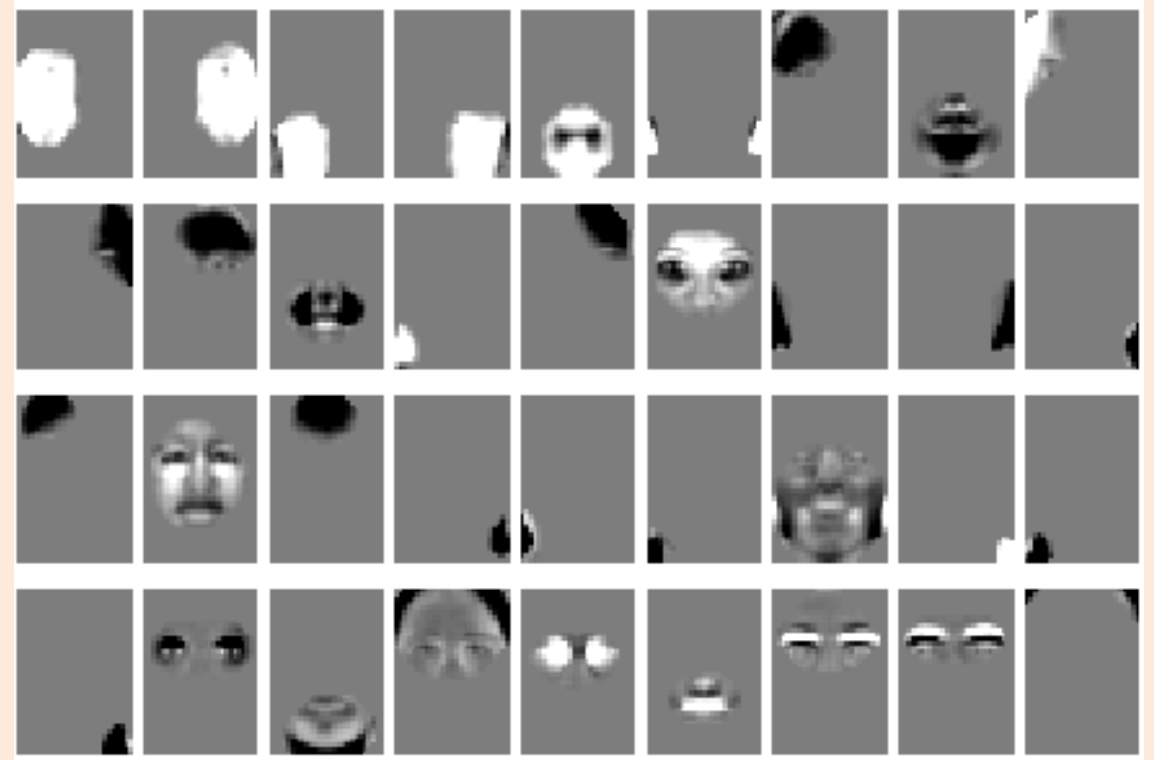PCA without orthogonality

sparsity due to non-negativity

sparsity due to L1-regularization

# Recent Work: Structured Sparsity

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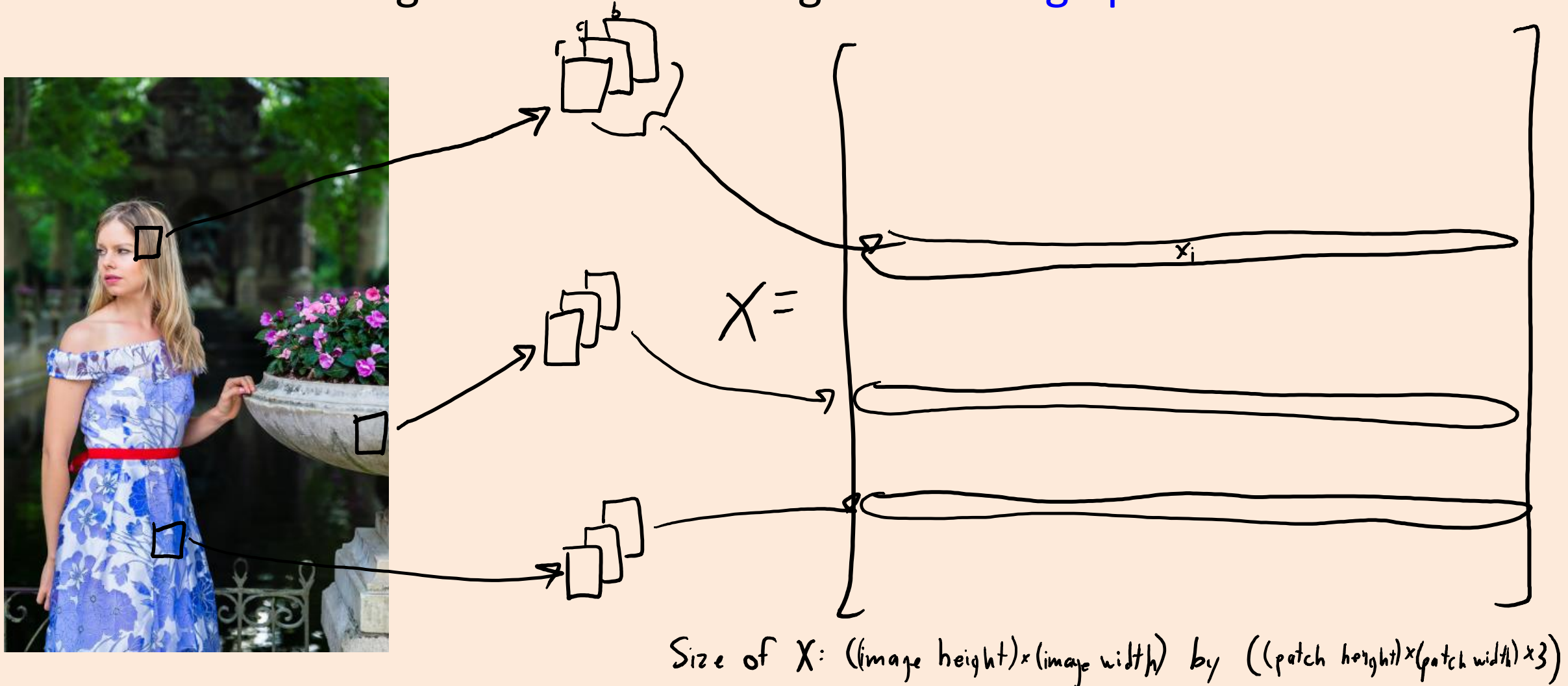


NMF

Sparse PCA with "structured" sparsity

# Summary

- Biological motivation for orthogonal and/or sparse latent factors.

- Non-negative matrix factorization leads to sparse LFM.

- Non-negativity constraints lead to sparse solution.
  - Projected gradient adds constraints to gradient descent.
  - Non-orthogonal LFMs make sense in many applications.

- L1-regularization leads to other sparse LFMs.

- Next time: the NetFlix challenge.

# Latent-Factor Models for Image Patches

- Consider building latent-factors for general image patches:



$X =$

Size of $X$: (image height) × (image width) by ((patch height) × (patch width) × 3)
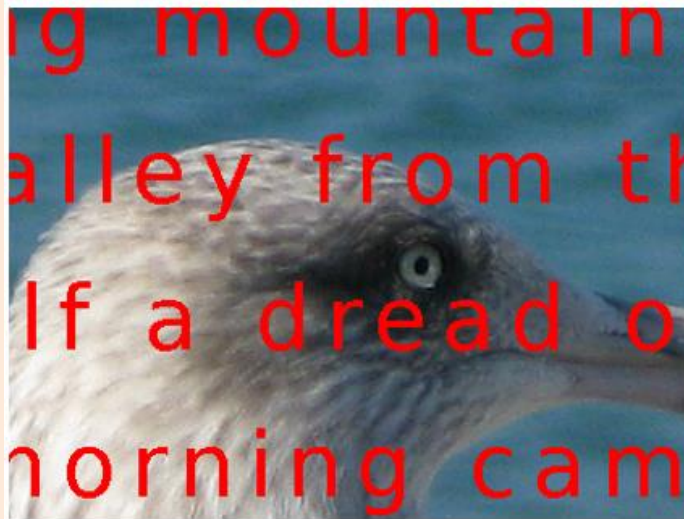
# Latent-Factor Models for Image Patches

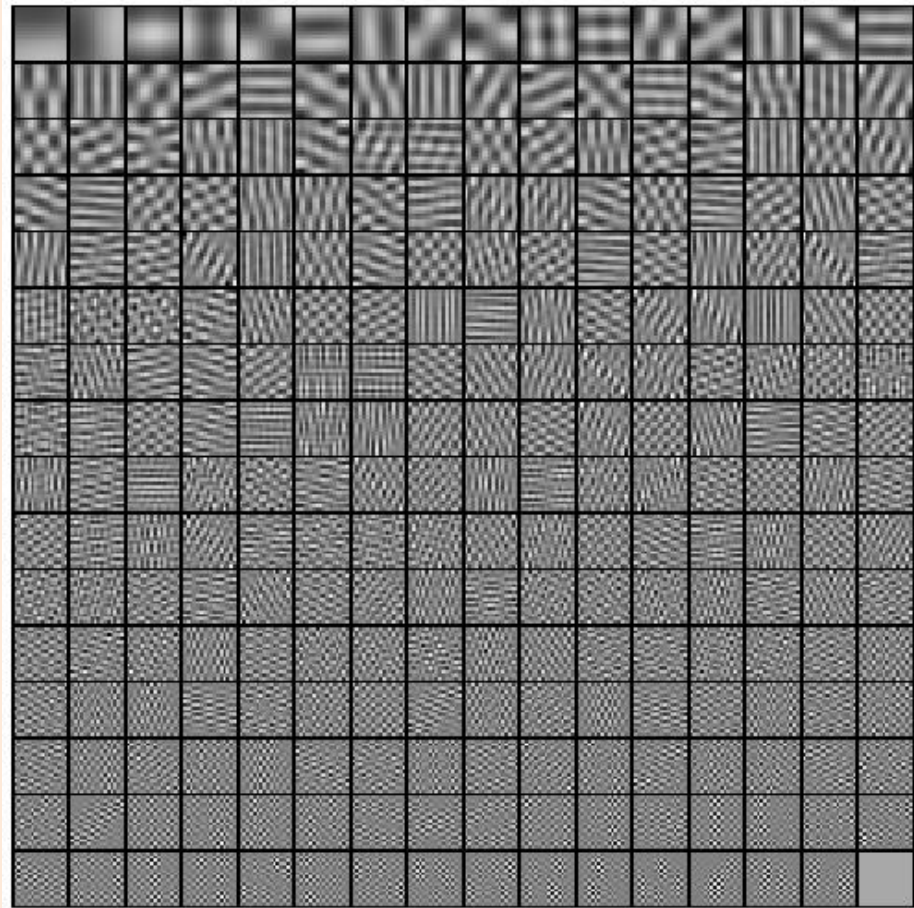- Consider building latent-factors for general image patches:



Typical pre-processing:

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

# Application: Image Restoration

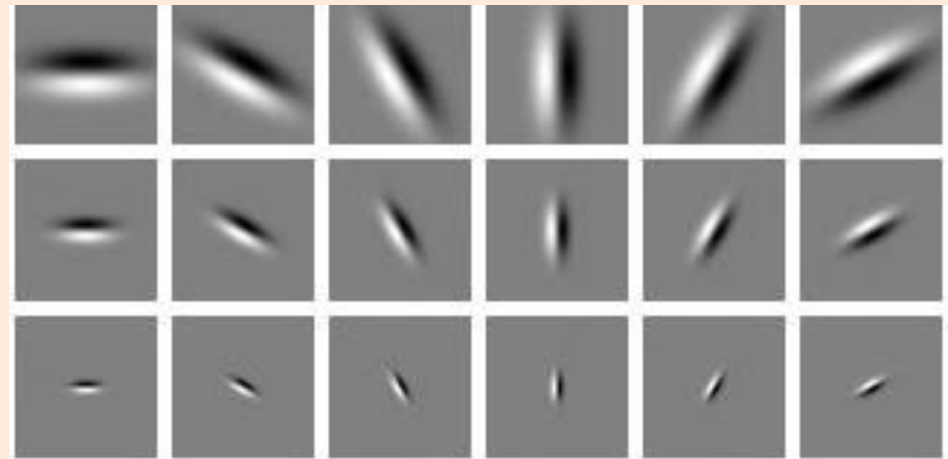# Latent-Factor Models for Image Patches



(b) Principal components.

Orthogonal bases don't seem right:
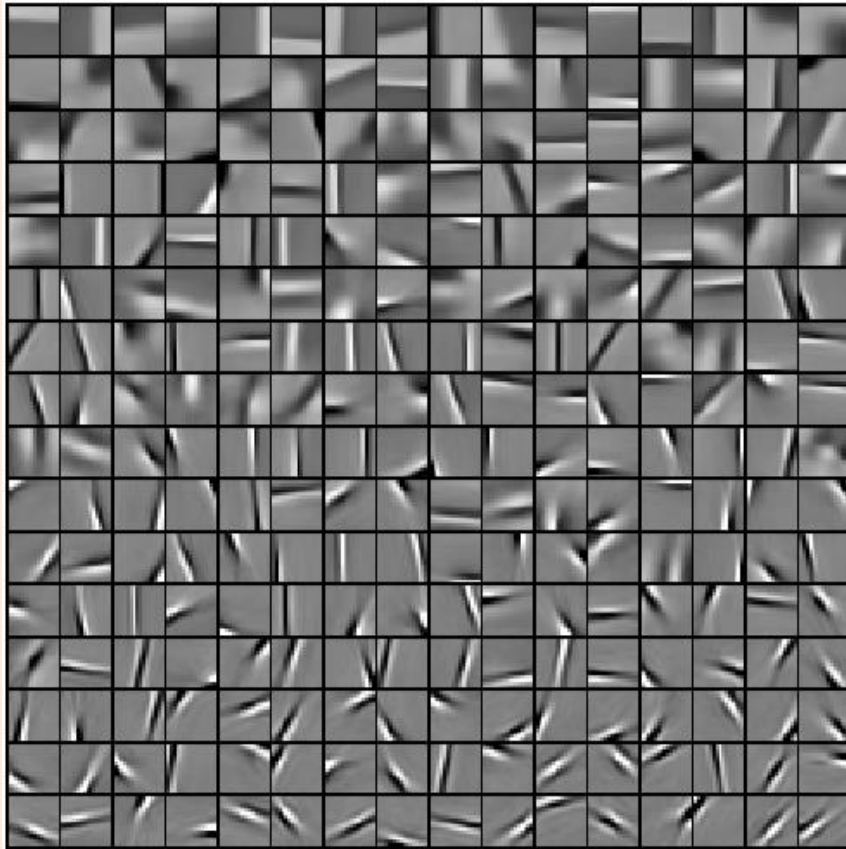- Few PCs do almost everything.
- Most PCs do almost nothing.

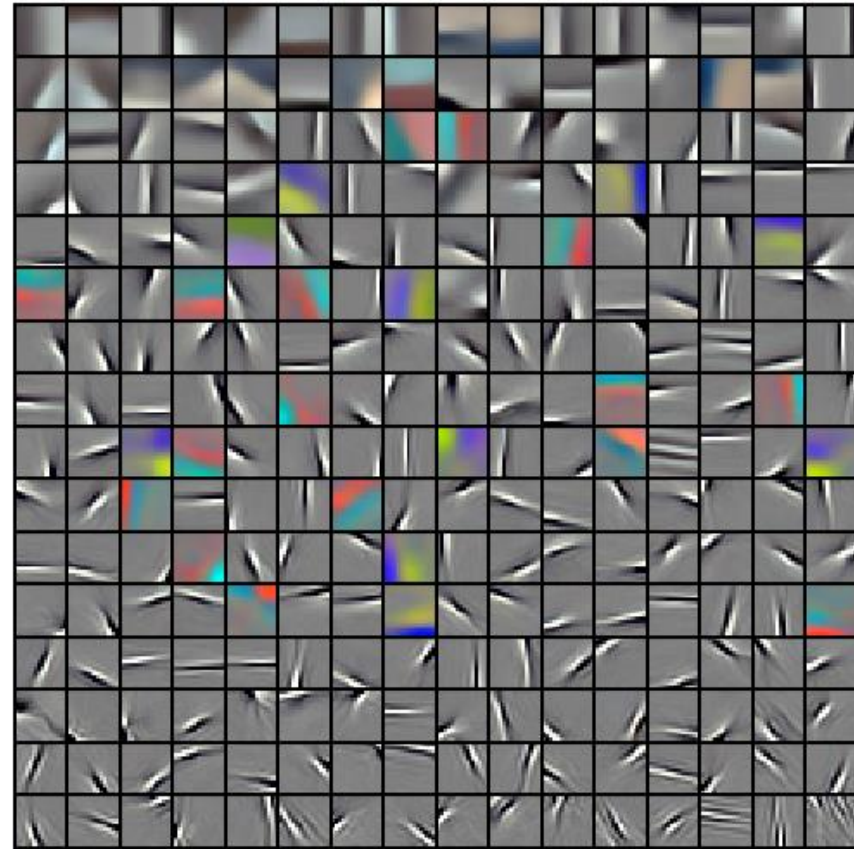We believe "simple cells" in visual cortex use:



'Gabor' filters

# Latent-Factor Models for Image Patches

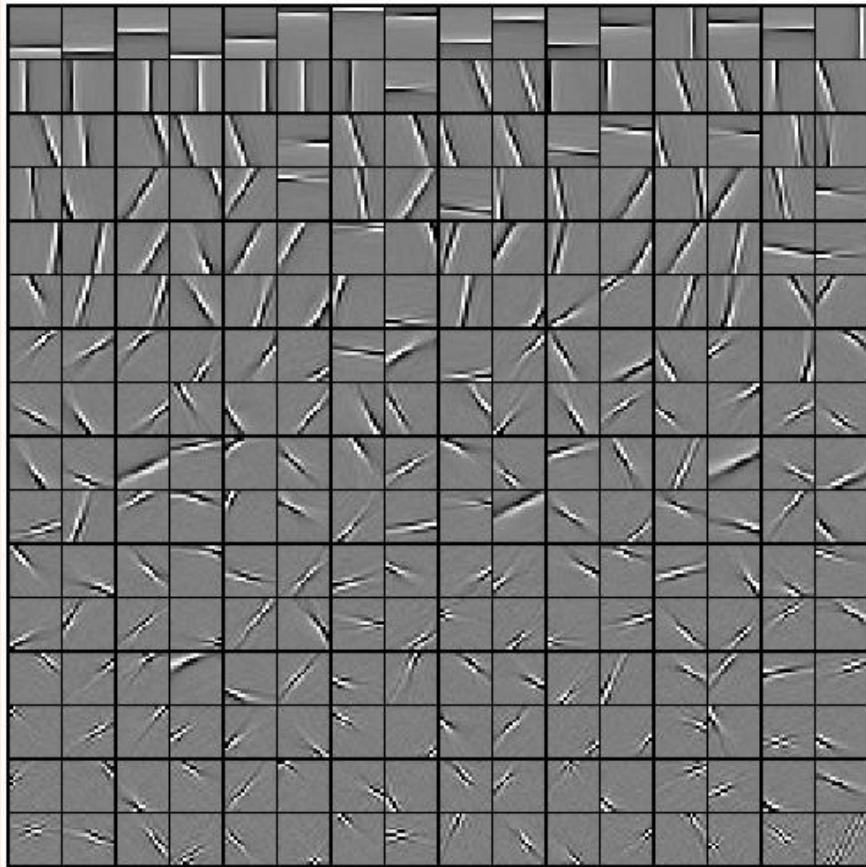- Results from a sparse (non-orthogonal) latent factor model:
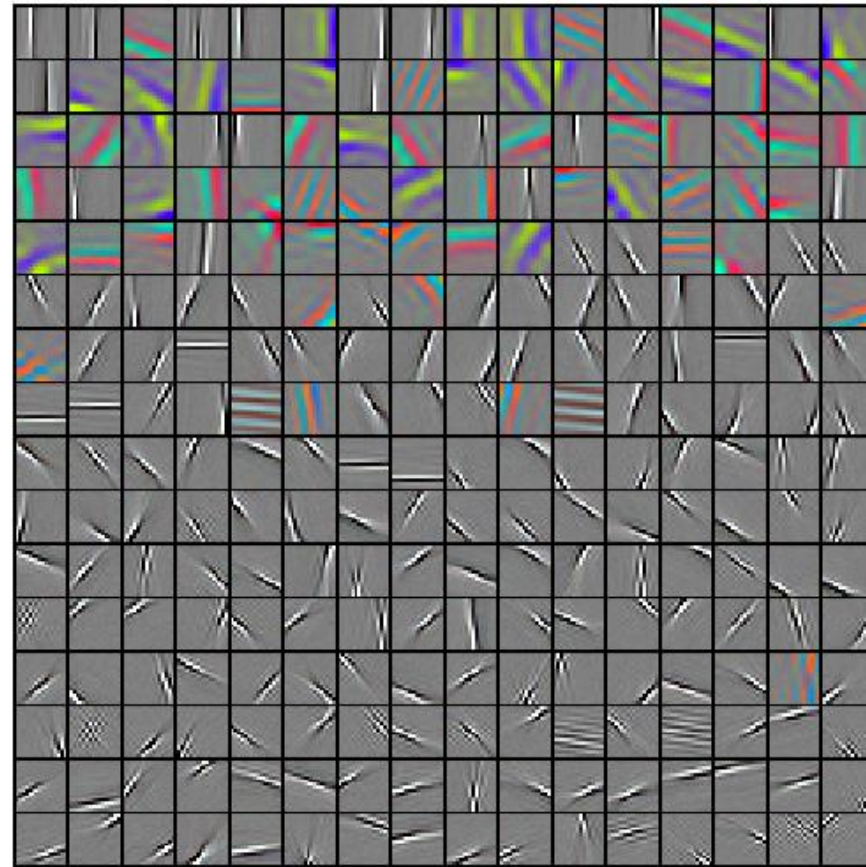


(a) With centering - gray.

(b) With centering - RGB.

# Latent-Factor Models for Image Patches

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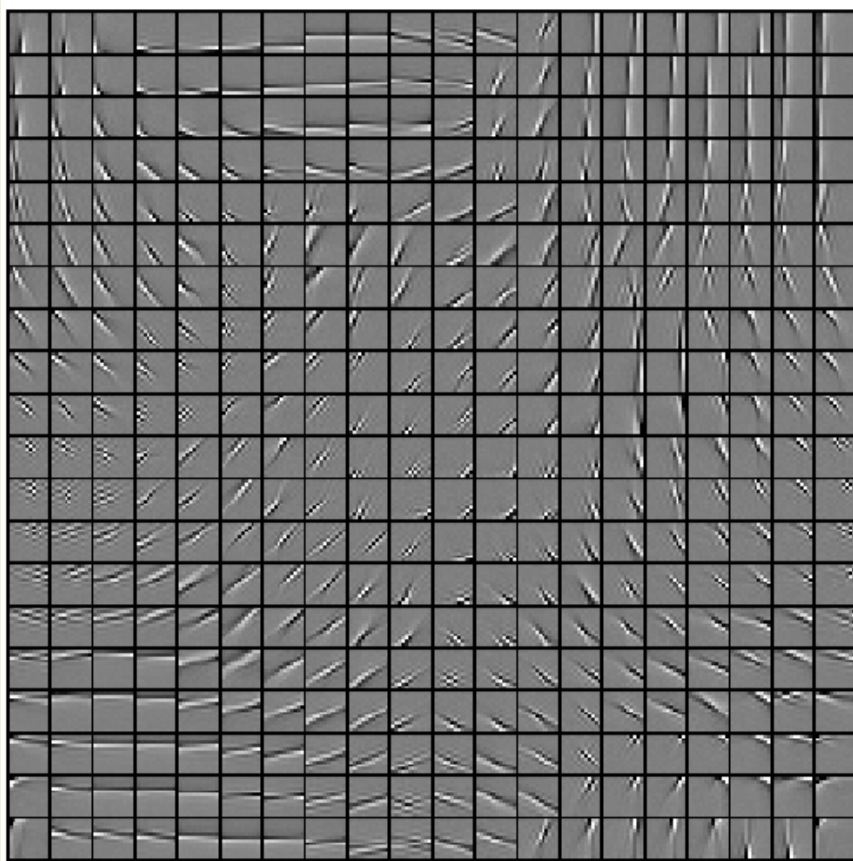


(c) With whitening - gray.

(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.