

# Image Data Compression

## Natural image statistics

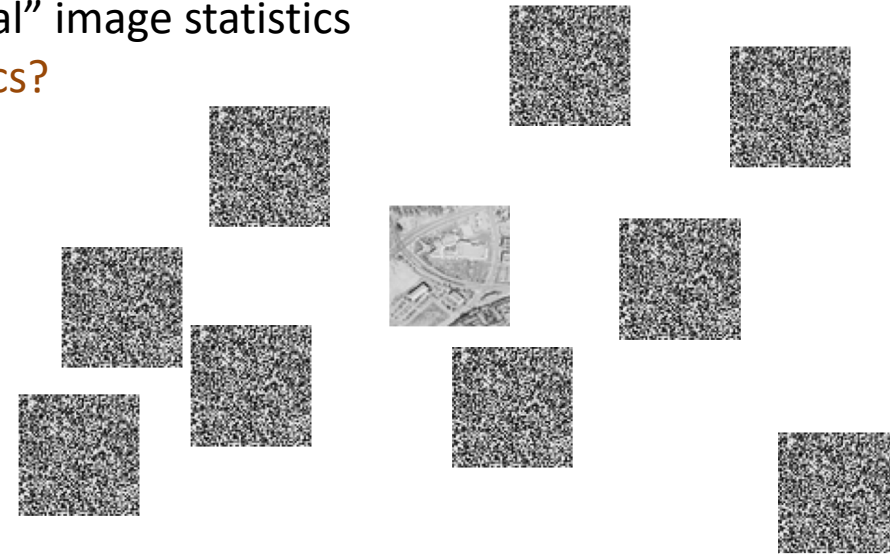
# Why do we need to study natural image statistics?

- Need good data models to have better compression methods
- Relatively recent field (the first comprehensive book on subject: 2009)
- Closely related to human visual system studies
  - (conjecture: biological visual systems are “optimal” for the “natural” images)
- Studies facilitated by big image databases and available processing power
- May better understand differences of “unnatural” image statistics
- Why is it challenging to study the image statistics?

Following the book:  
[Hyvärinen, Hurri, Hoyer  
Natural Image Statistics,  
Springer 2009]

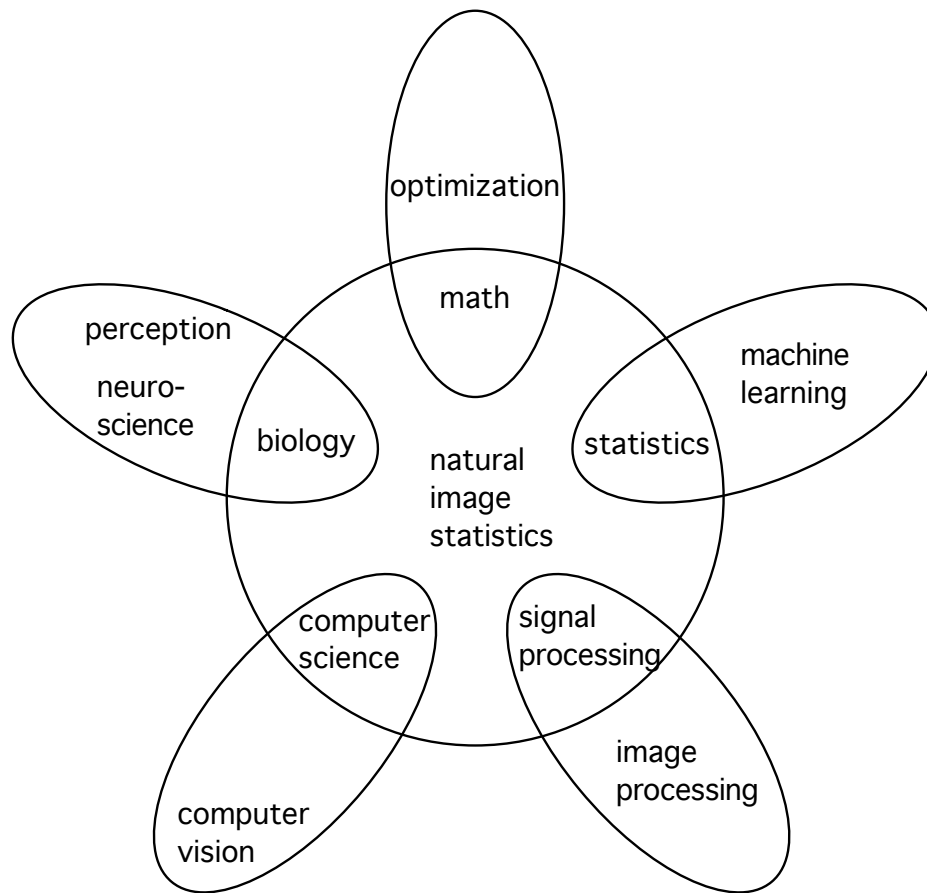
## A couple of (relatively) big numbers:

- Seconds passed since the Big Bang:  $\sim 10^{17}$
- Number of atoms in the Universe:  $\sim 10^{80}$
- Entropy of 65 x 65-pixel white noise images:  
 $H = 65 \times 65 \times 8 = 33800$  bits
- Number of 65 x 65 8-bit gray-scale images:  
 $2^H \approx 10^{10000}$
- Natural scenes are redundant, estimated entropy  $\sim 0.4 H$ ,
- => only one of  $10^{6000}$  white noise images has the basic statistics of natural scenes



*“The distribution of natural images is complicated. Perhaps it is something like beer foam, which is mostly empty but contains a thin mesh-work of fluid which fills the space and occupies almost no volume. The fluid region represents those images which are natural in character.”*

[Ruderman 1996]



## Models of natural image data

- Physical (generative) imaging models
- Non-linear manifold(s) of natural images
- Non-parametric sample-based models
- Biologically-inspired neural networks
- Simple models reproducing increasingly more complex statistics of natural images (equivalent to unsupervised learning)

**Goal: transform images to “convenient” spaces**

## Some proposed criteria of “good” models:

- “simplicity”  $\approx$  reduced dimensionality
- “sparseness”  $\approx$  compact representation
- “metabolic efficiency”  $\approx$  energetically efficiently processed by neurons in brain
- “learning efficiency”  $\approx$  easy to learn from unique images (each data sample is seen only once)
- “wiring length”  $\approx$  require minimal communication length between neurons

It all starts with statistics...

# Reminder of multivariate statistics (I)

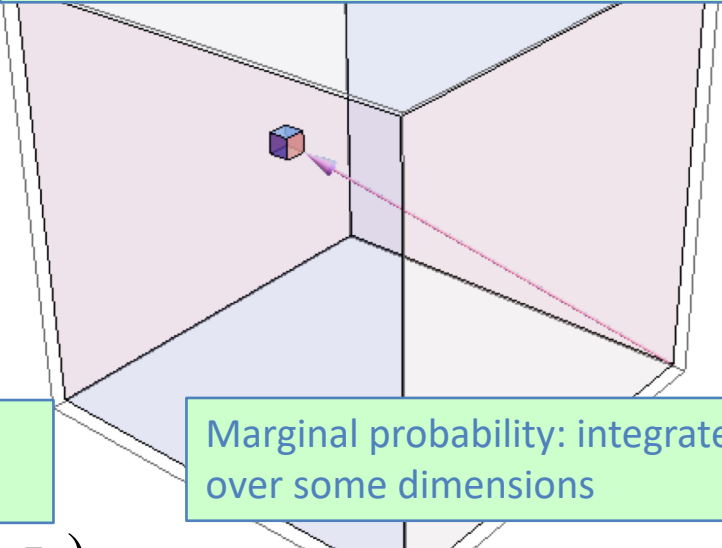
- Probability density function (PDF):

$$p_{\vec{z}}(\vec{a}) = \lim_{\Delta z \rightarrow 0} \frac{P(z_i \in [a_i, a_i + \Delta z] \forall i = 1, \dots, n)}{\Delta z^n}$$

$$\int p_{\vec{z}}(\vec{a}) d\vec{a} = 1$$

- Conditional vs marginal probability:

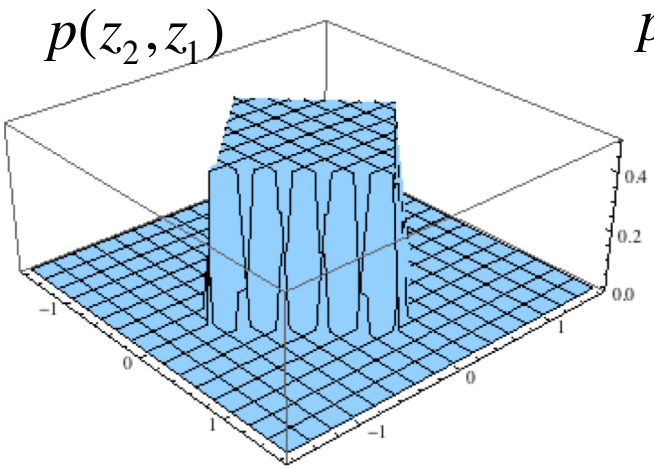
Slightly ambiguous, but shorter notation:  
 $p(\vec{z}) \equiv p(z_1, z_2, \dots, z_n) \equiv p_{\vec{z}}(\vec{z})$



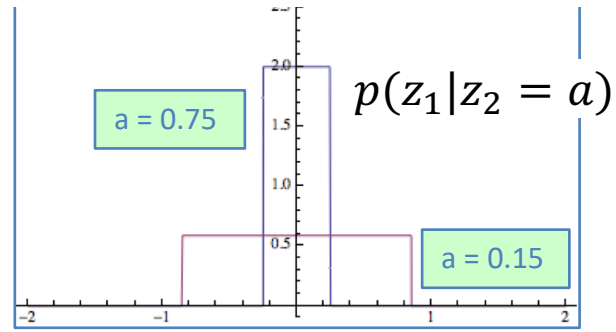
Joint probability:  
multi-dimensional PDF

Conditional probability:  
fix some dimensions

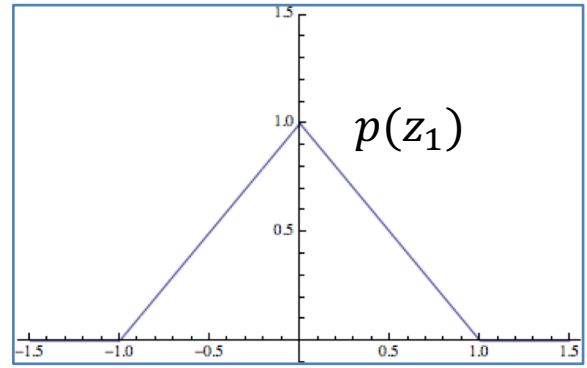
Marginal probability: integrate  
over some dimensions



$$p(z_2 | z_1 = a) = \frac{p(a, z_2)}{\int p(a, z_2) dz_2}$$



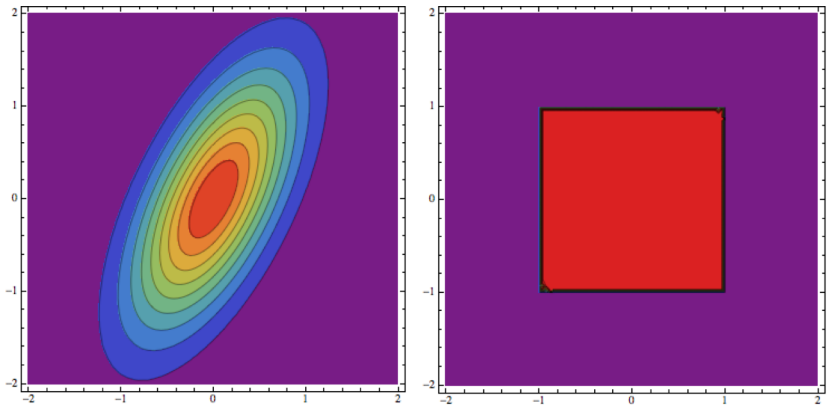
$$p(z_1) = \int p(z_1, z_2) dz_2$$



# Reminder of multivariate statistics (II)

- Variables  $z_1$  and  $z_2$  are independent iff their joint PDF factorizes:

$$p_{\vec{z}}(\vec{z}) \equiv p(z_1, z_2) = p(z_1)p(z_2)$$



- Expectation of a random vector (or its function):

$$E\{\vec{z}\} \equiv \langle \vec{z} \rangle = \int p_{\vec{z}}(\vec{a}) d\vec{a}$$

$$E\{g(\vec{z})\} = \int g(\vec{a}) p_{\vec{z}}(\vec{a}) d\vec{a}$$

1-st order statistic

**Note:** independent variables are uncorrelated, i.e.  
 $E\{g_1(\vec{z}_1)g_2(\vec{z}_2)\} = E\{g_1(\vec{z}_1)\}E\{g_2(\vec{z}_2)\}$   
 but uncorrelated variables may still be dependent!!!

- Variance and covariance (1-dimensional case):

$$\text{var}\{z\} = E\{(z - \langle z \rangle)^2\}$$

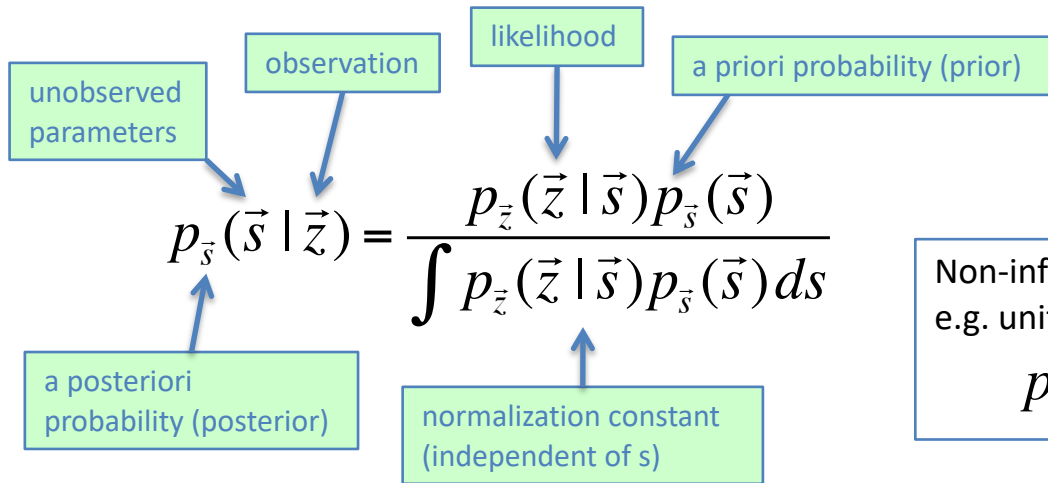
$$\text{COV}\{z_1, z_2\} = E\{z_1 z_2\} - \langle z_1 \rangle \langle z_2 \rangle$$

2-nd order statistics

- Covariance matrix (ND):

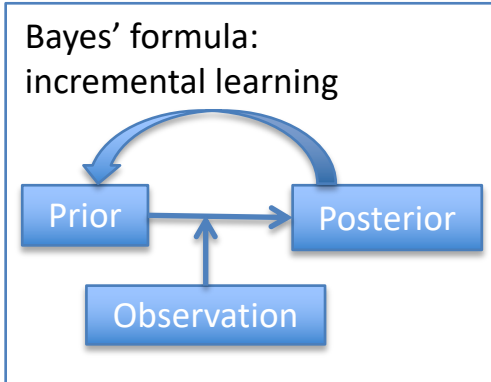
$$C(\vec{z}) = \begin{pmatrix} \text{COV}\{z_1, z_1\} & \text{COV}\{z_1, z_2\} & \dots & \text{COV}\{z_1, z_n\} \\ \text{COV}\{z_2, z_1\} & \text{COV}\{z_2, z_2\} & \dots & \dots \\ \dots & \dots & \dots & \dots \\ \text{COV}\{z_n, z_1\} & \dots & \dots & \text{COV}\{z_n, z_n\} \end{pmatrix} = E\{\vec{z}\vec{z}^T\} - \langle \vec{z} \rangle \langle \vec{z} \rangle^T$$

# Reminder of multivariate statistics (III)



$$p_{\vec{s}}(\vec{s} | \vec{z}) = \frac{p_{\vec{z}}(\vec{z} | \vec{s}) p_{\vec{s}}(\vec{s})}{\int p_{\vec{z}}(\vec{z} | \vec{s}) p_{\vec{s}}(\vec{s}) ds}$$

Non-informative prior:  
e.g. uniform distribution,  
 $p_{\vec{s}}(\vec{s}) = c$



## Parameter estimation (usually performed in log space)

- Bayes' formula:  $\log p_{\alpha}(\alpha | z_1, z_2, \dots, z_m) = \log p(z_1, z_2, \dots, z_m | \alpha) + \log p(\alpha) - \log p(z_1, z_2, \dots, z_m)$
- Maximum a posteriori probability (MAP):  $\alpha_{MAP}^* = \operatorname{argmax}_{\alpha} [\log p(z_1, z_2, \dots, z_m | \alpha) + \log p(\alpha)]$
- Maximum likelihood (ML):  $\alpha_{ML}^* = \operatorname{argmax}_{\alpha} [\log p(z_1, z_2, \dots, z_m | \alpha)]$  Constant, independent of  $\alpha$

Equivalent to MAP with constant prior

## Example: 1D Gaussian

$$p_z(z | \alpha) = \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{(z - \alpha)^2}{2}\right]$$

- Observe:  $z_1, z_2, \dots, z_m$ , iid:  $p_z(z_1, z_2, \dots, z_m | \alpha) = p(z_1 | \alpha) \cdot \dots \cdot p(z_m | \alpha)$
- Task: estimate  $\alpha$  (assume flat prior)

Likelihood:  $\log p(z_1, \dots, z_m | \alpha) = -\frac{1}{2} \sum (z_i - \alpha)^2 + const$

Solution via least squares method:

$$\alpha_{ML}^* = \frac{1}{m} \sum z_i$$

# Simple statistics of linear features

- One patch (e.g., 32 x 32 pixels) - random vector  $\vec{z}$  ( $n = 1024$  components)
- Patch at a random location in a random image from DB = one observation

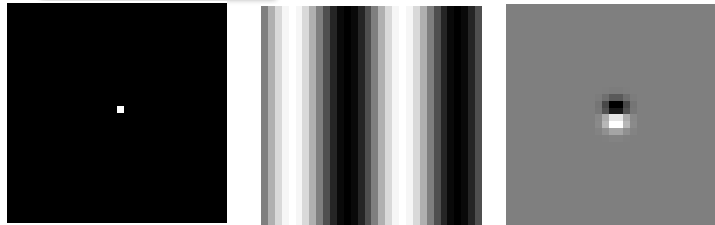
“Vectorize” image;  
gray-scale values

- Apply discrete linear transform:

$$s_i = \vec{w}_i^T \vec{z}$$

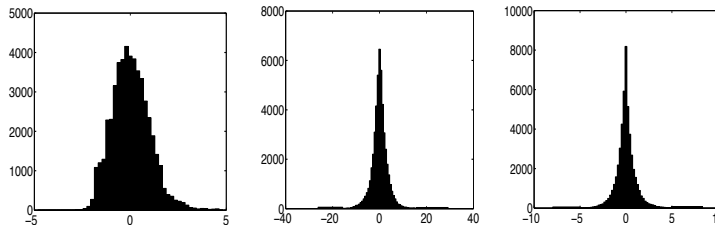
Basis functions, “features”, can also be interpreted as images

- Examples of features:



- Dirac filter (pixel value)
- Grating detector (cosine)
- Gabor edge detector

- Histograms of filter outputs:



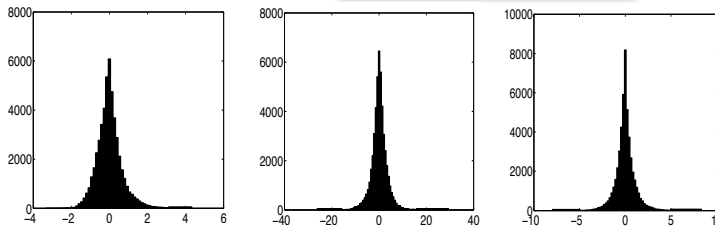
Two last histograms  
very non-Gaussian!

- Simplest image statistic: mean (DC component)  
(not to be confused with the mean over sample!)

$$\hat{\vec{z}} = \vec{z} - \frac{1}{n} \sum_{i=1}^n z_i$$

The only first-order structure

- Mean-subtracted histograms:



Grating and edge detectors  
are orthogonal to mean value,  
Dirac filter is not orthogonal

In what follows, always assume subtracted mean

General strategy: remove parts that are well-understood, and study what remains

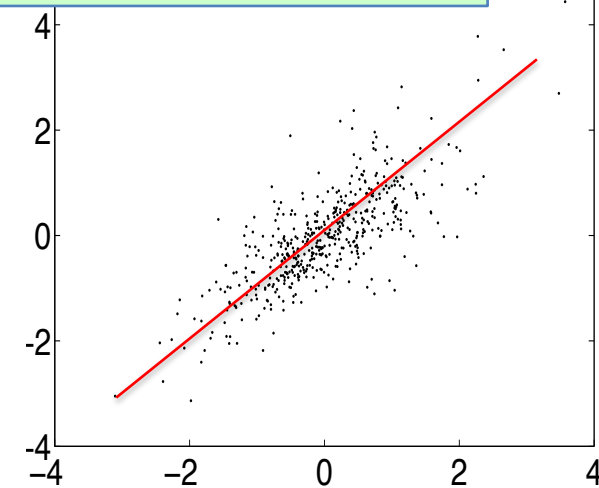
# Principal component analysis

- Same as KLT applied to whole n-dimensional patches
- Second-order structure (pixel-to-pixel correlations)
- Principal components “capture” the largest data variances
- Technical requirements: mean-subtracted images, limit norm of  $w$  by e.g. 1

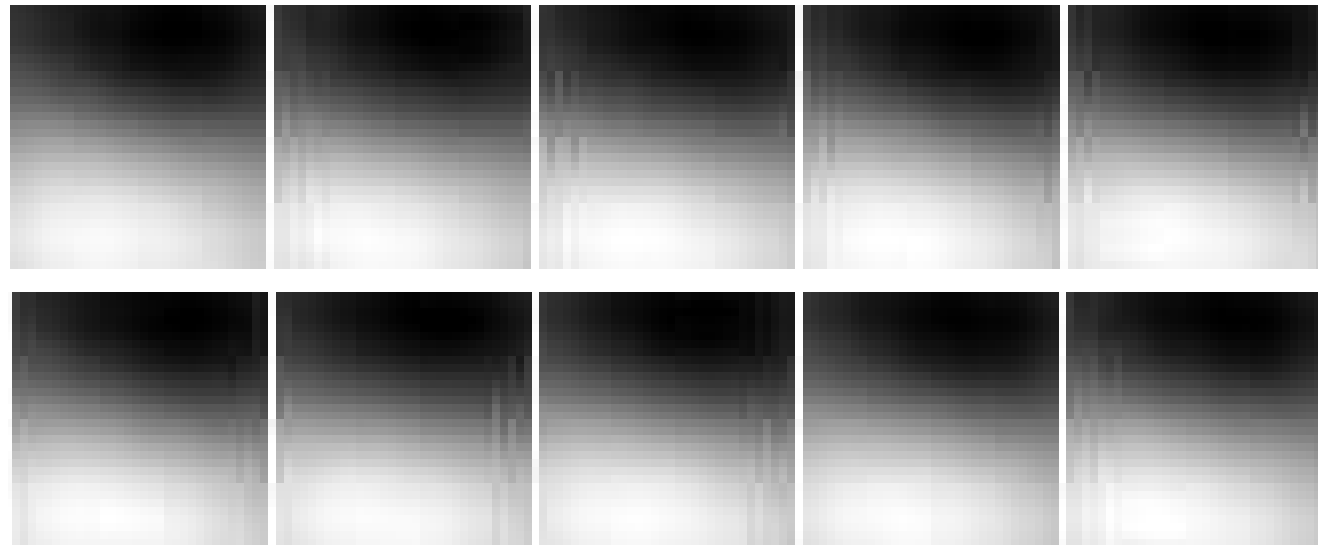
**Learning the first principal component:** max output variance in a given sample,

$$\vec{w}^* = \operatorname{argmax}_{\vec{w}, \|\vec{w}\|_2=1} \left[ \frac{1}{T} \sum_{t=1}^T (\vec{w}^T \vec{z})^2 \right]$$

Scatter plot of two adjacent pixel values and their primary axis



First principal component computed for 32x32 patches for 10 different randomly sampled datasets:



Seems to be rather stable, but not extremely useful per se...

# Principal component analysis (PCA) – further components

**Deflation-type definition:** all following vectors maximize variance in the orthogonal subspace to the previous ones

- Resulting vectors are orthogonal
- Resulting projections are uncorrelated
- Equiv. to frequency separation!

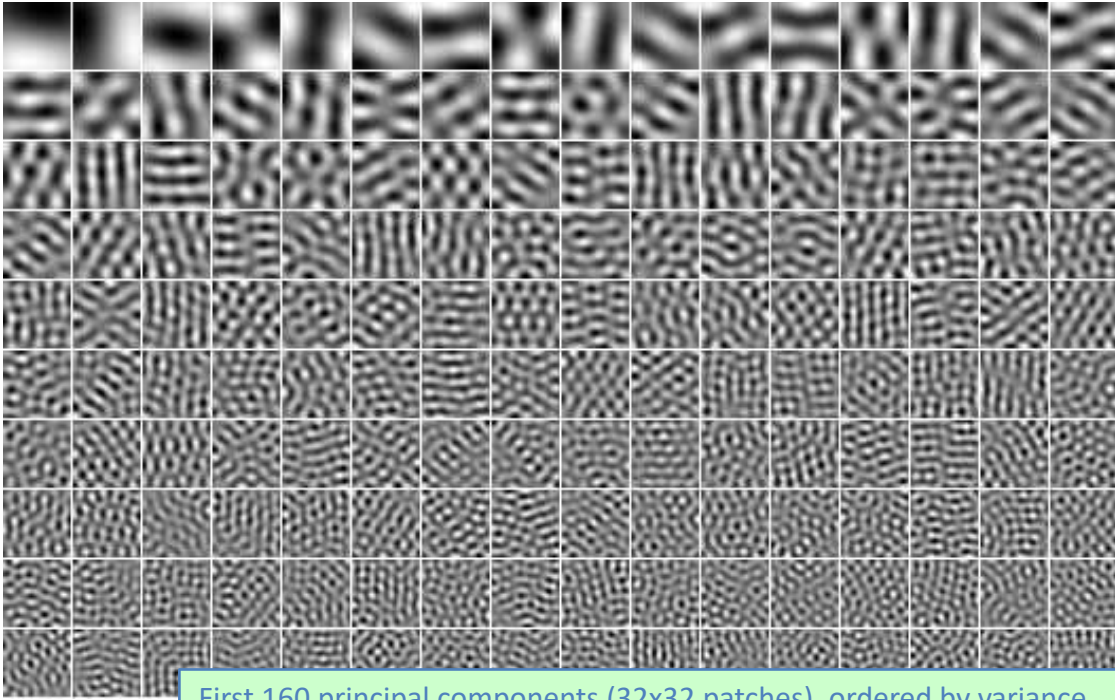
$$\vec{w}_j^* = \operatorname{argmax}_{\substack{\vec{w}, \|\vec{w}\|_2=1, \\ \vec{w} \perp \vec{w}_1^*, \dots, \vec{w}_{j-1}^*}} \left[ \frac{1}{T} \sum_{t=1}^T (\vec{w}^T \vec{z}_t)^2 \right]$$

### Practical calculation:

- Eigenvectors of covariance matrix (Eigenvalues give variance)
- Cf. Karhunen-Loève transform
- Alt: Fourier of correlation matrix

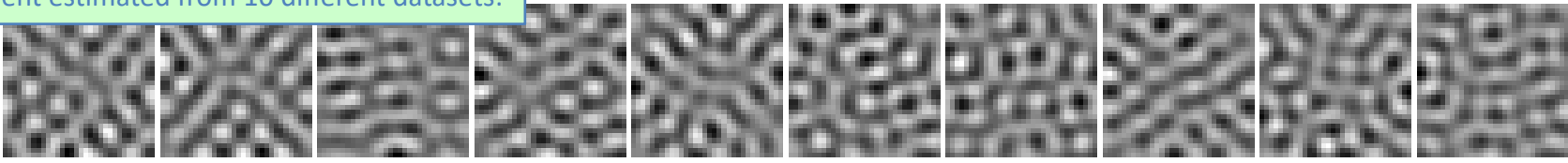
### Problems with PCA:

- No reasonable property of images
- Similar variances for higher-order components -> unstable!
- Depends on random fluctuations



First 160 principal components (32x32 patches), ordered by variance

100<sup>th</sup> component estimated from 10 different datasets:



# Why is PCA useful (I)?

## Dimension reduction with PCA

- All pixel values ( $n = 32 \times 32 = 1024$ ) - too many variables!
- Apply a linear transform, reduce to  $m < n$  values:

$$\vec{z} = (z_1, \dots, z_n)^T, \quad s_i = \vec{w}_i^T \vec{z}, \quad z_j = \vec{a}_j^T \vec{s}, \quad i = 1, \dots, m, \quad m < n;$$

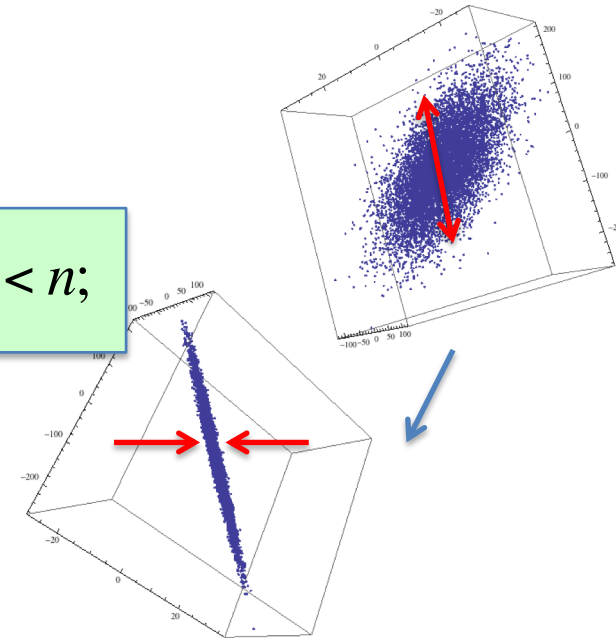
- New variables must “preserve information”; in other words:
- In reconstructing  $z$ , minimize average squared error:

$$E \left\{ \left\| \vec{z} - \sum_i \vec{a}_i \cdot s_i \right\|^2 \right\} \rightarrow \min$$

- Assumptions: orthonormal transformation basis,

$$\vec{w}_i^T \vec{w}_j = \delta_{ij}$$

- **Solution:** take  $m$  first principal components!
- Defined up to arbitrary rotation in  $m$ -dim space
- Non-unique definition of directions, but more stable principal subspace (spanned by principal directions)



## Proof for the 1<sup>st</sup> principal component:

$$s = \vec{w}^T \vec{z}, \Delta = E \left\{ \left\| \vec{z} - s \cdot \vec{w} \right\|^2 \right\}$$

$$\Delta = E \left\{ \left\| \vec{z} \right\|^2 \right\} + E \left\{ (\vec{w}^T \vec{z}) \cdot (\vec{w}^T \vec{z}) \cdot \left\| \vec{w} \right\|^2 \right\} - 2E \left\{ (\vec{w}^T \vec{z}) \cdot (\vec{w}^T \vec{z}) \right\}$$

$$\left\| \vec{w} \right\|^2 = \vec{w}^T \vec{w} = 1,$$

$$\Delta = E \left\{ \left\| \vec{z} \right\|^2 \right\} - \vec{w}^T E \left\{ \vec{z} \cdot \vec{z}^T \right\} \vec{w} = \text{var}(\vec{z}) - \vec{w}^T C(\vec{z}) \vec{w}$$

min  $\Delta$  :

$$L = \vec{w}^T C(\vec{z}) \vec{w} - \lambda \cdot (\vec{w}^T \vec{w} - 1),$$

$$\frac{dL}{d\vec{w}} = 2 \cdot \vec{w}^T C(\vec{z}) - 2\lambda \cdot \vec{w}^T = 0, \Leftrightarrow [C(\vec{z}) - \lambda I] \vec{w} = 0.$$

For natural images: ~ 10% of dimensions are enough!

# Why is PCA useful (II)?

## Whitening by PCA – generic steps

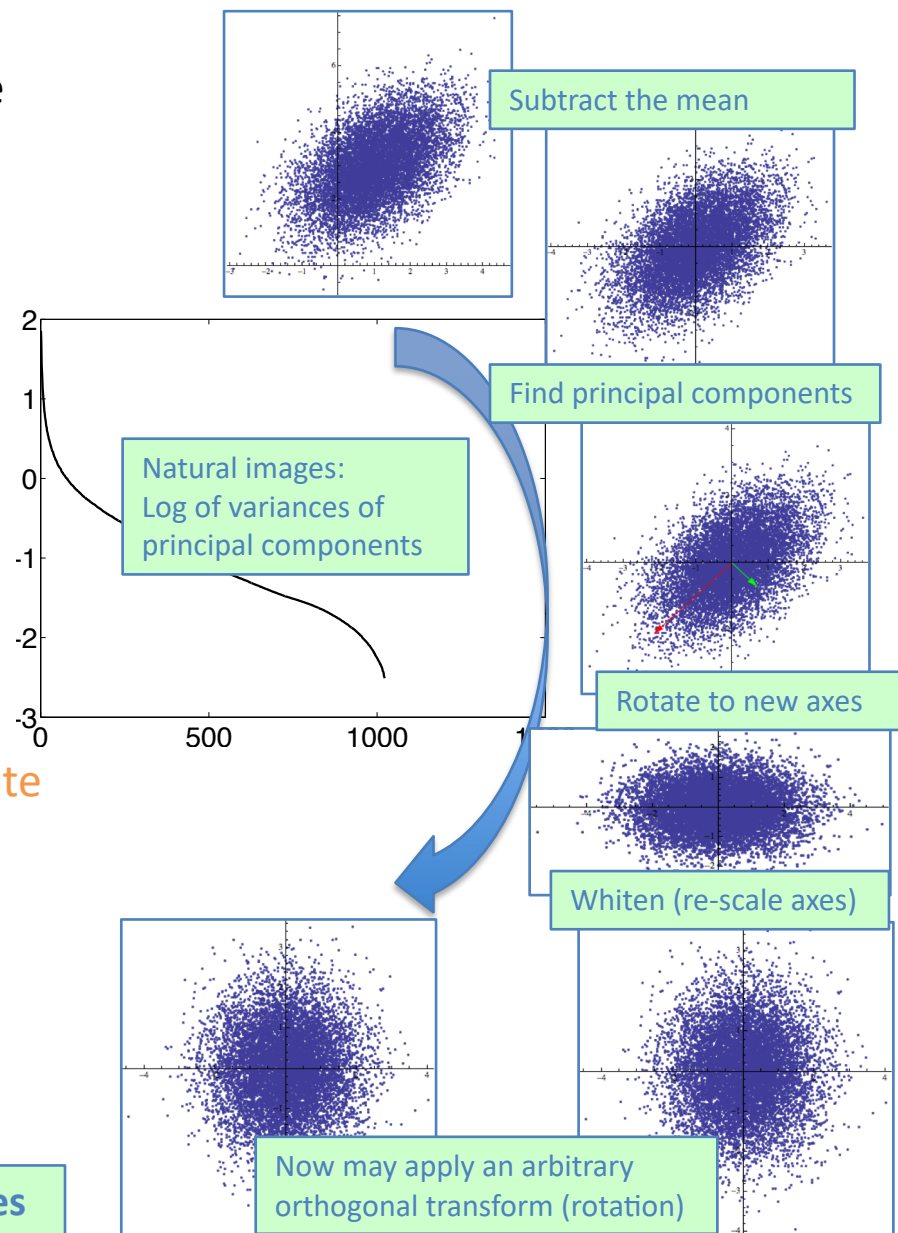
- Goal: transform an image to variables which are uncorrelated and have normalized variance:

$$E \{ s_i s_j \} = \delta_{ij}$$

- Solution: re-scale the principal components,

$$y_i = \frac{s_i}{\sqrt{\text{var}(s_i)}}$$

- Completely exploits and removes second-order information (correlations and variances)
- Standard pre-processing of statistical data
- There exist many whitening transforms; in fact, any orthogonal rotation of whitened data is white
- White distribution with the highest entropy – Gaussian:  $\vec{y} \sim N(\vec{0}, I)$
- Zero correlation of features = orthogonality of their vectors in the whitened space

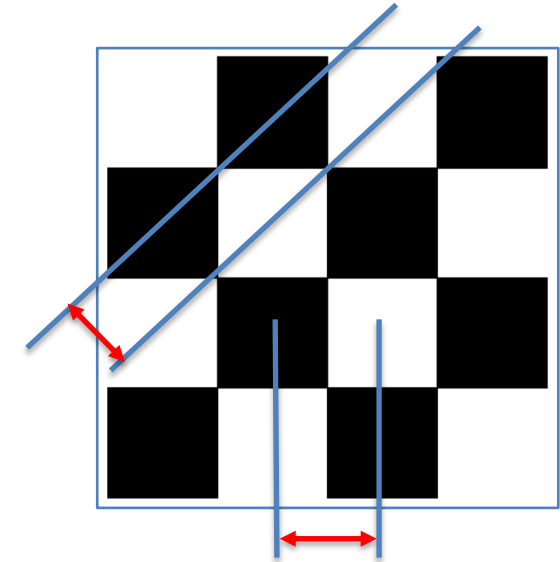


Further analysis: study non-Gaussianity of natural images

# Why is PCA useful (III)?

## Anti-aliasing of rasterized images by PCA:

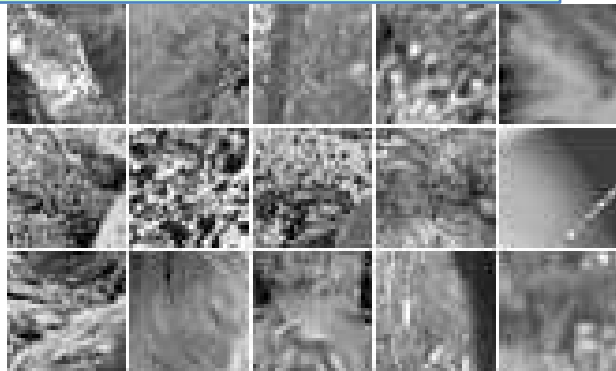
- Highest H or V frequencies can have only two phases
- Along diagonals, max frequency is  $\sqrt{2} = 1.414\dots$  times higher
- Non-isotropic representation of isotropic natural images!
- Simple dimensional reduction by PCA simultaneously filters oblique high frequencies away, improves anisotropy
- PCA can also be formulated in Fourier space



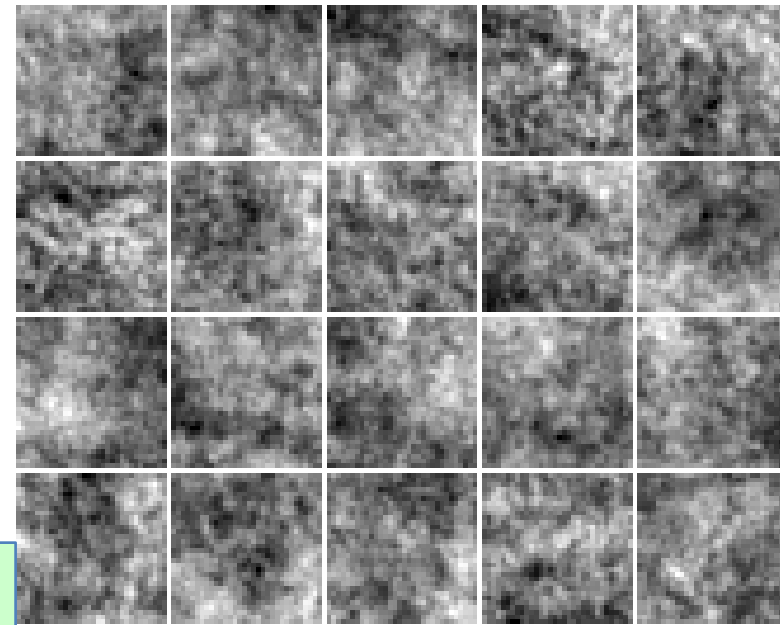
## PCA as a generative model:

- Idea: generate random components from Gaussian distribution with  $s_i$  according to estimated variances
- Perform inverse transform, i.e. “hallucinate” images
- Verdict: relatively poor generative quality!

Natural randomly sampled 32x32 images



Random 32x32 images generated from PCA model



# Beyond second-order statistics: sparse features

## Problems with PCA:

- Features do not resemble neural receptive fields
- Generative model not extremely successful

## Need higher-order properties. E.g.: sparseness

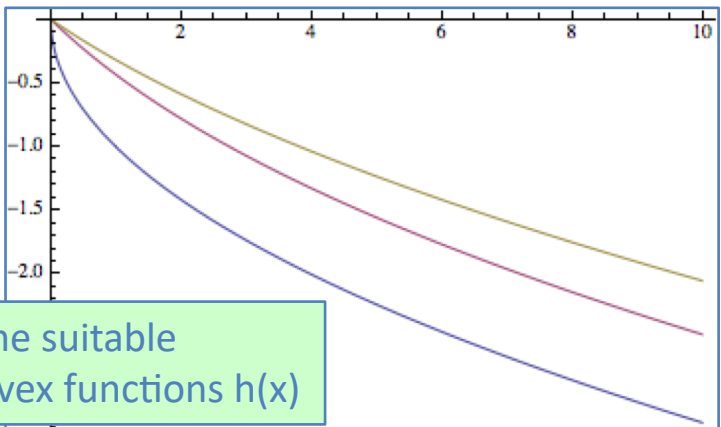
- Approximate definition: random variable components most of time are small (“rarely active”)
- **Not the same as small variance!**

## May use many measures of sparseness

- Maximize some non-linear function of squared variable,

$$E \{ h(s^2) \} \rightarrow \max$$

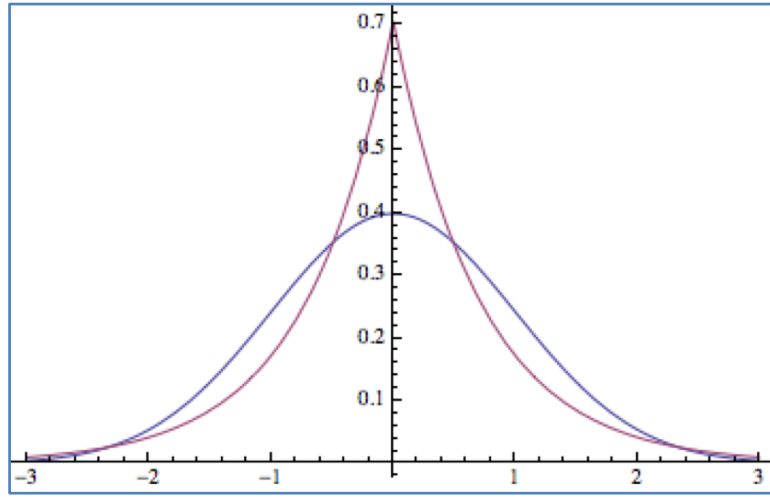
Need  $h(x)$  to be **convex**: near zero – higher weights



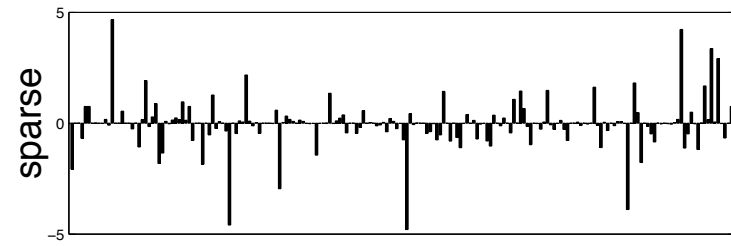
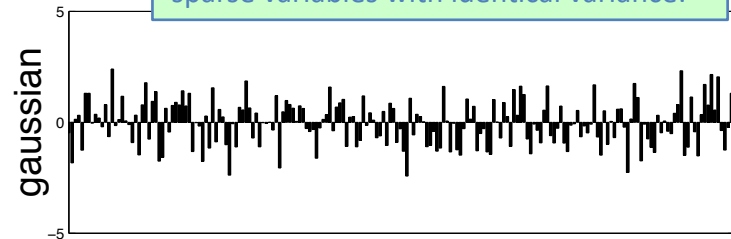
Some suitable convex functions  $h(x)$

- Commonly used:**
- kurtosis (4-th moment),
  - $-x^{1/2}$
  - $-\log \cosh(x^{1/2})$
  - $-(x + e)^{1/2}$
- Optimal measure:**
- $\log p(x^{1/2})$

A Gaussian and some sparse PDFs with the same variance



Samples of normally distributed and sparse variables with identical variance:



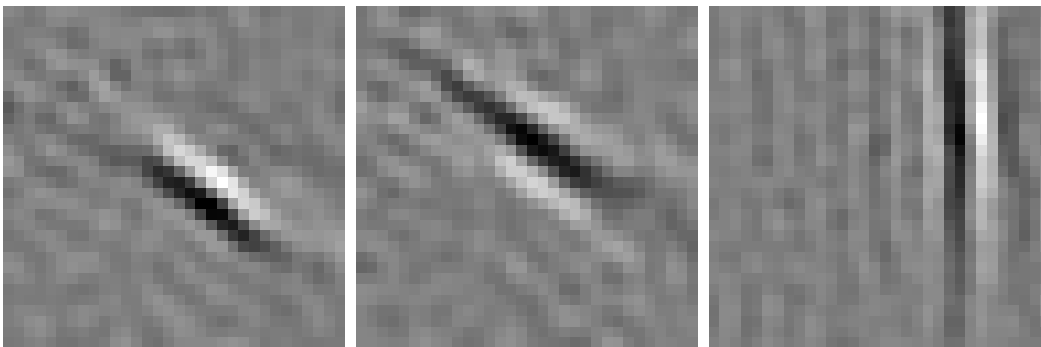
# Linear feature detectors with maximum sparseness

- Start from PCA-preprocessed (whitened) data
- Choose some convex function  $\mathbf{h}(\mathbf{x})$
- Learn one feature that maximizes sparseness:

$$\vec{v}^* = \arg \max_{\vec{v}, \|\vec{v}\|=1} E \left\{ h \left( (\vec{v}^T \vec{s})^2 \right) \right\}$$

## Results (finding a single sparse feature):

- Features localized in space, frequency, orientation
- Resemble receptive fields of neurons
- Problem: many local maxima, all well-localized!



Weights found by sparseness maximization in natural images, different initialization (i.e. local maxima)

## Learning multiple sparse features

- Start from different points (can find the same max many times!)

## A better method: “deflation”

- Find local maxima satisfying some constraints
- One possible choice: subsequent features should be un-correlated with the previously found ones  
(= orthogonal in the whitened space),

$$\vec{v}_j^* = \operatorname{argmax}_{\vec{v}, \|\vec{v}\|=1} E \left\{ h \left( (\vec{v}^T \vec{s})^2 \right) \right\},$$
$$E \left\{ (\vec{z}^T \vec{v}_j^*) (\vec{z}^T \vec{v}_i^*) \right\} = 0, \forall 1 \leq i < j$$

- Leads to gradual deterioration of features (too strict constraints);
- Last vectors have too little space to optimize

# Symmetric de-correlation

**Better yet:** maximize sum of individual sparseness measures, under constraints of unit variance and symmetric de-correlation:

$$\{\vec{v}_1, \dots, \vec{v}_n\}^* = \operatorname{argmax}_{\{\vec{v}_1, \dots, \vec{v}_n\}} \sum_{i=1}^n E \left\{ h \left( (\vec{v}_i^T \vec{s})^2 \right) \right\}$$
$$E \left\{ (\vec{v}_i^T \vec{s})(\vec{v}_j^T \vec{s}) \right\} = \delta_{ij}$$

Number of features is limited by the pre-processing step (PCA)

**A small philosophical problem:**

**Above: Sparseness of feature** = only few feature values over sample images (as a function of t)

**Really want: Sparseness of representation** = feature values over index i for a single image

- Similar to spoken language: many words, each phrase contains only a few of them
- Number of features can exceed dimensionality of data!

Under these conditions,

$$\sum_{i=1}^n h \left( (\vec{v}_i^T \vec{s})^2 \right)$$

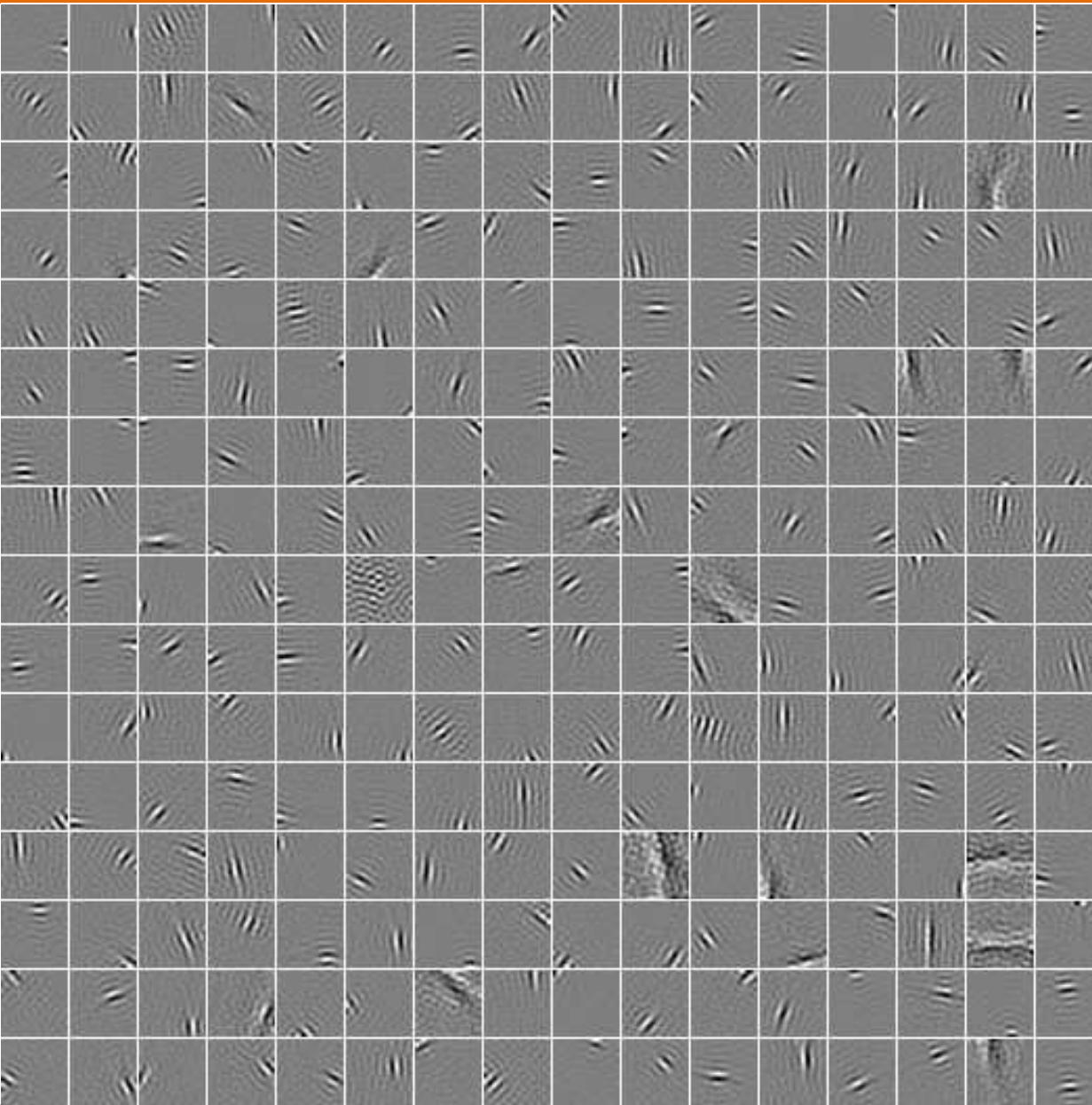
measures the sparseness of a single image, sum over sample equivalent to the formula above

Both definitions equivalent if following conditions hold for a typical image:

- The mean of features is approximately zero
- Mean of square of features equals approximately one

**This holds if features are statistically independent and have identical distribution (iid)**

# Sparse coding feature detectors from natural images



## Computation:

- DB of 50000 32x32 patches
- No ordering (symmetric de-correlation)
- PCA reduction to  $m = 256$
- FastICA algorithm to learn sparse features

## Features:

- Localized in space
- Well-defined orientation
- Multi-scale (small and large)
- Similar to neural receptive fields, with distributions over frequency, phase, ...
- Many entries are just shifted copies of other features

Is that the best we can do?

# Why is sparseness useful?

## So far:

- A better statistical model of input data
- Efficient coding of images (many zeros after transform, easy to compress)
- Sparse features resemble receptive fields of simple cells in a biological visual system
- Some plausible explanation: neurons try to minimize their firing rates to save energy

## Still have open questions:

- Sparseness measure was chosen ad hoc. What is the optimal function  $h(x)$ ?
- Why was de-correlation (of many features) needed?
- What is the true Bayesian prior distribution of images?

## Further refinement: generative models

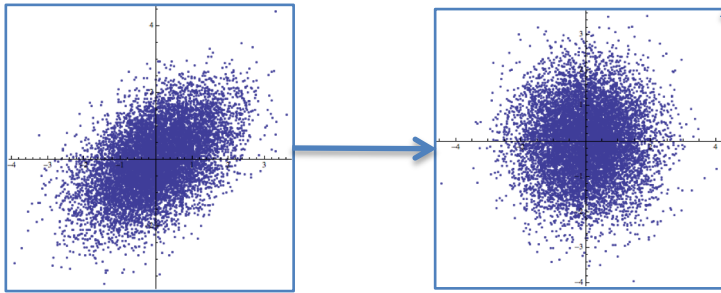
- Observed data (images) generated as transformations of hidden (latent) variables
- Specific flavor of model: Independent Component Analysis (ICA)
- ICA in some specific formulation is equivalent to finding maximally sparse features!

**Main idea:** assume that transformed coefficients are statistically independent over sample,

$$s_i = \vec{w}_i^T \vec{z}, \quad p(s_1, \dots, s_n) = p_1(s_1) \cdot \dots \cdot p_n(s_n), \quad \text{var}(s_i) = 1$$

# Why does PCA not produce independent components?

## Gaussian distribution: PCA and whitening

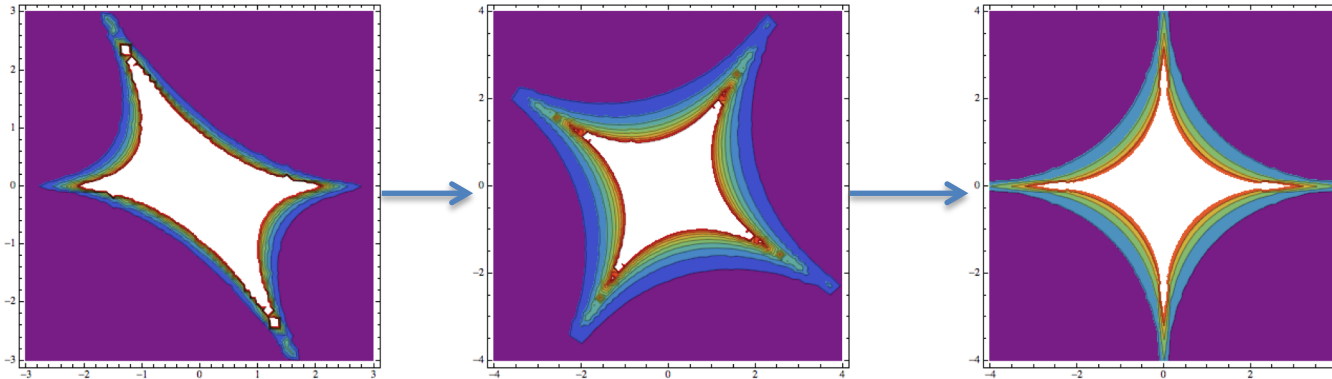


- All rotations are equivalent
- The PDF spherically symmetric
- Uncorrelated Gaussian variables are already independent!

### Using only second-order information:

- The correlation matrix is symmetric, i.e.  $n(n+1)/2$  values
- A transform matrix needs  $n^2$  independent values
- Remaining DoFs:  $n(n-1)/2$  rotation angles

## Sufficiently non-gaussian distribution: PCA and whitening, followed by ICA:



### Benefits of ICA:

- Strong relation to sparseness
- Provides optimal sparseness measure
- Justifies de-correlation
- PDF for Bayesian inference

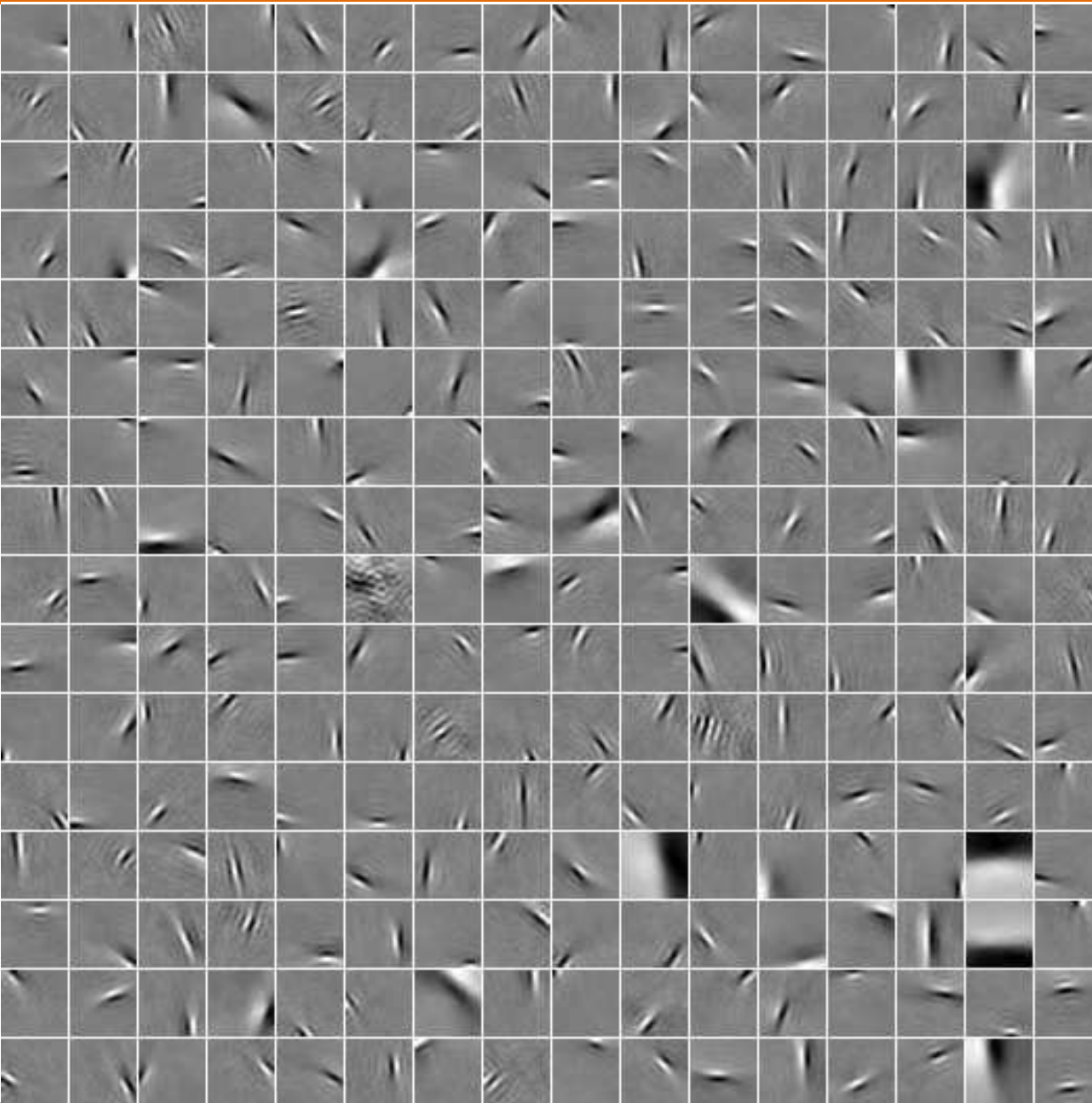
## ICA feature learning: maximum (log) likelihood estimation

$$s_i = \vec{v}_i^T \vec{z}, p(s_1, \dots, s_n) = p_1(s_1) \cdot \dots \cdot p_n(s_n), \Rightarrow p(\vec{z}) = |\det(V)| \cdot \prod_{i=1}^n p_i(\vec{v}_i^T \vec{z})$$

$$\{\vec{v}_1, \dots, \vec{v}_n\}^* = \arg \max_{\vec{v}_1, \dots, \vec{v}_n} \left[ \log(|\det(V)|) + \sum_i \log(p_i(\vec{v}_i^T \vec{z})) \right]$$

Many methods exist to find this maximum!

# ICA-based features learned from natural images



## Computation:

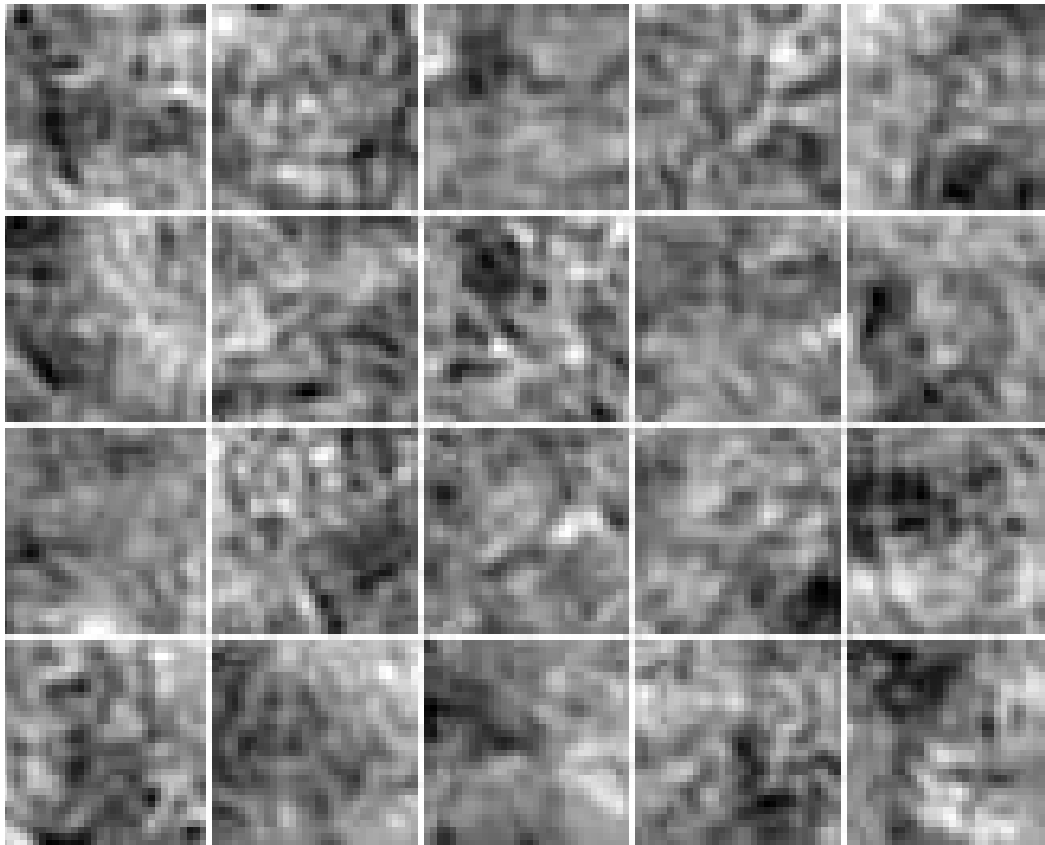
- DB of 50000 32x32 patches
- No ordering (symmetric DC)
- PCA reduction to **n = 256**

## Features:

- Similar to sparse features
- Same properties wrt localization as before
- Frequency channels seem to be statistically independent!
- Extracted features still exhibit some dependence!

Use as a generative model:  
sample “typical” natural images

# Image patches generated from the ICA model



## Computation:

- Marginal distributions over each feature estimated from real images

## Results:

- Better than PCA: edge-like structures
- Still, not extremely “natural” 😊

## Further (current) research directions:

- Minimum-entropy coding
- Over-complete bases
- Non-negative models
- Non-linear features, energy detectors (cf. complex cells in visual cortex)
- Independent Subspace Analysis (ISA)
- Multi-layer models
- Modeling extra-striate cortex (V2, ...)
- Markov Random Field models
- ...

Still an open field, more work is required!