

Computational Vision

U. Minn. Psy 5036

Daniel Kersten

Lecture 24: Object Recognition

Initialize

Spell check off

```
In[1]:= Off[General::spell1];
```

Outline

Today

Object recognition

Considerable accumulation of knowledge from computer science, neuroscience, and behavioral studies over the past 20 years. Object recognition in computer vision remains a top-priority project, in part fueled by practical uses for internet image searches. In neuroscience, there have been extensive studies of the role of the primate ventral visual stream in object processing, and the proposed “end-goal” of basic or “core” recognition.

We will focus on the computational problems, and consequences for models of human recognition. And the role of geometric modeling in theories of object recognition. Discussion of ideal observer analysis for 3D object recognition

High-level vision, visual tasks

Let’s set the context for object recognition and review of levels of abstraction & task-dependency.

Intermediate-level vision

Low-level vision is what can be inferred from image intensity patterns without knowledge that the patterns are caused by surfaces, their properties and relationships.

Intermediate-level vision refers to inferences possible without knowledge of an object’s category, using

only generic object and surface properties. The processes include:

Selection of features likely to have a common cause

Generic, organizational processes, for surface grouping, involving Gestalt principles

Occlusion (“domain overlap”) & relative surface depth

Perceptual integration

 Cue integration (weak fusion)

 Cooperative computation (strong fusion)

Attention

Intermediate-level processes are useful for interpreting novel objects and scenes. These processes could also be useful for feature extraction to be used to store in memory, and later for testing match against stored representations. The idea is that more abstract features and relations would be more robust to image variations, of the sort discussed below.

High-level vision

Functional tasks

 Object recognition--familiar objects

 entry-level, subordinate-level

 Object-object relations

 Scene recognition

 Spatial layout

 Viewer-object relations

 Object manipulation

 reach & grasp

 Heading, time-to-contact

Task dependency: explicit (primary) and generic (secondary, nuisance) variables

One can't think of invariance without considering what has to be discounted for a given type of task to achieve it. Consider several classes of scene causes of image pattern I.

 Image = f(shape, material, articulation, viewpoint, relative position, illumination, foreground and background clutter)

Which variables are more important to estimate precisely for various tasks?

Task: Object Recognition (labelling)

 I=f(**shape, material**, articulation, viewpoint, relative position, illumination, foreground and background clutter)

Distinguish: detection, categorization, and identification. We take a closer look below.

Task: Absolute depth (e.g. for reaching)

$I=f(\text{shape, material, articulation, viewpoint, position, illumination, foreground and background clutter})$

Task: grasp

$I=f(\text{shape, material, articulation, viewpoint, relative position, illumination, foreground and background clutter})$

Problem: all the scene variables contribute to the variations in the image

We will focus on shape-based-recognition

Shape-based object recognition:

- estimate geometrical shape (primary variables)

- discount sources of image variation not having to do with shape (secondary variables)

e.g. integrating out geometrical variables such as translation, rotation, and scale, but also photometric variables such as illumination, to estimate shape for object recognition.

We'll postpone detailed discussion of the crucial problems that variations due to background clutter, and within-category shape also need to be taken into account.

Object recognition: computational issues

Analysis of image variation

Which variables are important to estimate for recognition depends on the level of abstraction required.

Variation within a subordinate-level category

What distinguishes a mallard from a wood duck? Honeycrisp from a Braeburn? Doberman from an Alsation?

Think of all the ways the images (or “appearances”) of an object, like a **male mallard** duck might vary. List the generative causes.

- illumination

- level, direction, source arrangement, shadows, spectral content

- view

- scale

- translation

- 2D & 3D rotation

- articulation

- non-rigid,

- e.g. joints, hinges, facial expression, hair, cloth, wings,
- physical size
 - small and big apples, shoes, dogs, ...
- background (segmentation)
 - bounding contour (affected by variation in pattern of intensities over the boundary regions)
- occlusion (segmentation)

What is left?

Geometric, metric properties of shape are important for subordinate

- e.g. sensitivity to "configurational information",
- but also material (e.g. orange vs. lemon)

Prototype representations -> what kind of model for variation?

Problem: Given training on only a discrete set of views of an object, how does vision generalize to other views of the same object?

Variation within a basic-level category

What distinguishes any duck from any dog, or from a chair or an apple?

Shape is important -- but qualitative, rather than metric aspects important.

Part types and their spatial relationships is one answer. E.g. geons and geon relations (Biederman).

Distinct prototypes each with a space of deformations is another approach.

How to achieve part-relation invariance within a category? I.e. most coffee cups have a hollow cylinder and a handle. While the

shapes of parts and relative positions can vary somewhat, there is a level of invariance that constrains an instance of a cup to be labeled as a "cup".

We will look at a structural relations theory, and diagnostic "fragment" theory

Fragment-based methods or "features of intermediate complexity": Ullman, S., Vidal-Naquet, M., & Sali, E. (2002). Learning informative features for class categorization.

Variation within a super-ordinate category

Variation is primarily semantic, and non-pictorial, rather than perceptual. Examples include: bird, mammal, furniture, fruit.

Categories can also be defined by affordances, e.g. the space of objects that one can sit on.

Levels of abstraction in human recognition

Basic-level vs. subordinate-level

How is the distinction relevant to human behavior?

Behavioral experiments (Rosch et al.)

Issue of prototypes with a model for variation vs. parts and relations.

e.g. average image face, the most familiar

visual priming methods to tease apart distinct representations (e.g. Biederman)

Neuropsychological (Damasio and Damasio)

temporal lobe lesions disrupt object recognition

fine-grain distinctions more easily disrupted than coarse-grain ones

e.g. Boswell patient—can't recognize faces of family, friends, unique objects, or unique places. Can assign names like face, house, car, appropriately.

Also superordinate categories: "tool"

prosopagnosics

faces vs. subordinate-level. Categorization but not identification.

neural evidence for distinction? IT hypercolumns?

Other factors that affect recognition

Computational complexity

E.g. a crumpled piece of paper -- difficult to compute and remember exact shape.

A set of diagnostic features, or summary statistics, may be more relevant.

Object ensemble

E.g. imagine only one red thing in your world--no need to process its shape. You may get the correct basic and subordinate-level categories.

Context is important

Small red thing flying past the trees.

High "cue validity" for male Cardinal bird

Context can even over-ride local cues for identity

Sinha and Poggio, Nature 1996



Democrat coalition.



See too: Cox, D., Meyers, E., & Sinha, P. (2004). Contextually evoked object-specific responses in human visual cortex. *Science*, 304(5667), 115-117.

Getting a good image representation

For object recognition, the contributions due to the secondary or "generic variables", (e.g. illumination and viewpoint) need to be discounted, and depending on the level of abstraction, various object features such as shape and material may need to be estimated. How?

-- Measurements of image information likely to belong to the object. This principle should constrain segmentation.

regions with similar textures, super-pixels (Shi and Malik, 2000; Sharon et al., 2006)

problems with: specularities, cast shadows, attached shadows (from shading).

edge detection is really noisy, and ambiguous as to cause, so what are these image "features"?

although noisy, are edges/groupings sufficiently reliable to determine object class?

-- Cue integration to improve estimates of where object boundaries are located:

combine stereo, motion, chromatic, luminance, etc..

-- Incorporate intermediate-level constraints to help to find object boundaries or "silhouettes".

Gestalt principles of perceptual organization (symmetry, similarity, proximity, closure, common fate, continuity,..)

long smooth lines (David & Zucker, 1989; Shashua & Ullman, 1988; Field and Hess, 1993)

-- "cooperative computation" for object shape, reflectance and lighting.

"intrinsic images" of Barrow and Tenenbaum

explaining away, e.g. for occlusion

A bottom-up procedure for perfect segmentation or edge-parsing is not known, but is also not needed for basic-level recognition. Accurate boundary segmentation may be important for subordinate-level decisions, such as identification of an object as being the same, or only slightly different.

Open questions regarding shape representations

Object geometry--Surfaces & shape, small scale surface structure

How can we describe objects themselves in terms of their geometry?

Are objects represented in terms of a view-dependent, dense local representation of shape, e.g. local surface orientation?

Or intrinsic properties, such as curvature?

Or in terms of parts? What is the relationship of parts of objects to each other?

Compositional representations

Role in object recognition, e.g. structural descriptions

To what extent can intermediate-level computations be short-circuited, if the task is narrower--e.g. category labeling with approximate localization?

We'll come back to this when we discuss neural-network inspired, feedforward models of recognition, trained using deep learning methods.

Storing information about an object and matching stored information to new appearances

In the next section, we take a closer look at several ways in which to store information about 3D objects in a way that is useful for recognizing novel instances or views of these previously seen objects.

Structural description: high-level features or parts plus relations. How would you describe the letters

A, ~~A~~, A, A, and A that is independent of font? Strokes and their spatial relationships.

Image-based: low-level features plus metric comparisons, transformations?

2D views?

3D object-centered?

Given a representation of the image information likely to be due to 3D object in memory, how does the brain store, then later when given another view, index and verify? Consider two extremes:

Nearest neighbor to 2D views?

Transformation of 3D model to fit 2D view?
Or something in between?

Two broad classes of models for object recognition

Image-based models

2D image description S , memory model M

-Matching

Try various M 's to test whether $S = F(M)$? Here one imagines the brain using its generative model F , to test for models of S that match S .

Or one could try to anticipate all the ways in which images of M might vary, and allow for those feedforward:

test: $F^{-1}(S) = M$? Here one imagines the brain has processes that could operate feedforward to get the input in the right invariant format to test against memories.

Current feedforward deep convolutional networks fall into this category of models, where the input is progressively transformed into a representation that can be tested and mapped onto a category label.

What are possible representations of M ? And how to model F ? Note that M could more or less resemble a 2D or 3D representation.

Image-based or "Exemplar" theories
view-specific features are stored in memory

Image-based models predict view-point dependence (e.g. Rock & DiVita (1987), and in general,

dependence on frequency of experience with particular views.

Poggio & Edelman, 1990; Bülthoff & Edelman, 1992; Tarr & Bülthoff, 1995; Liu, Knill & Kersten (1995); Troje & Kersten (1999)

In neural network models, the representations are highly constrained by the assumed neural architecture.

E.g. a hierarchy of spatial convolutions, with pooling and sigmoidal non-linearities.

A key issue is how efficiently can a recognition model learn new objects perhaps from only a few exposures, and recognize these objects later from completely new views or "appearances". The next class of models tries to address this issue.

Structural description and compositional models

These models emphasize the importance of explicitly representing spatial relationships. In psychology, see Hummel, J. E., & Biederman, I. (1992). Dynamic binding in a neural network for shape recognition.

Psychological Review, 99(3), 480–517. And in computer vision, see Geman, S., Potter, D. F., & Chi, Z. (2002). Composition systems. Quarterly of Applied Mathematics, 60(4), 707–736, and Ullman (1996).

Structural description theories use invariants to find parts (assumption is that this is easier than for the whole object), build up description of the relations between the parts which description specifies the object. E.g. a triangle shape, the letter "A" (three parts, with two "cross relations" and one "cotermination" relation".

Could be based on 2.5 D sketch => object-centered representation that is independent of viewpoint?
e.g. Marr's generalized cylinders

Some versions predict view-point independence. An well-known example is Biederman's geon theory.

Biederman (1987) proposed a basic vocabulary of parts from the extraction of invariants, "non-accidental properties", such as:

co-linearity of points or lines => colinearity in 3D

cotermination of lines=>cotermination in 3D (e.g. Y and arrow vertices)

skewed symmetry in 2d=>symmetry in 3D

curved line in 2D =>curved line in 3D

parallel curves in 2D => parallel in 3D (over small regions)

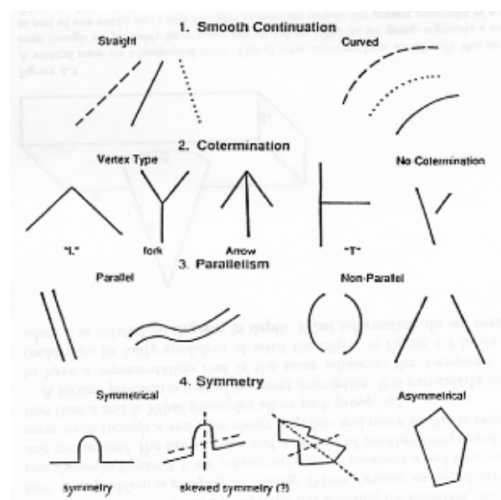
The considerations of non-accidental image properties lead to the idea of objects being represented in terms of elementary "parts" or

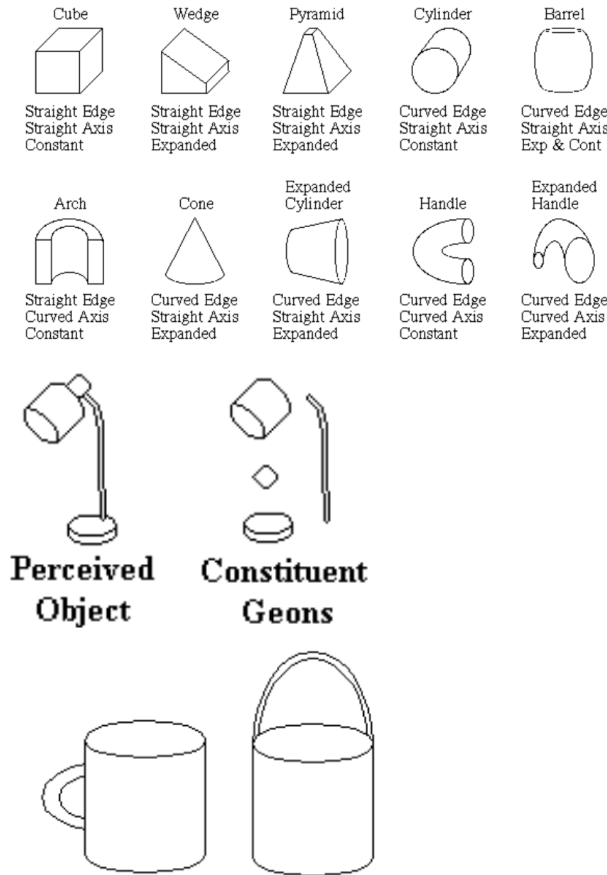
=> geons (box, cylinder, wedge, truncated cone, etc..) and a description of their spatial relationships to each other.

partial independence of viewpoint

Figures below from: Biederman, I. (1987). Recognition-by-components: A theory of human image understanding.

Psychological Review, 94, 115-147.



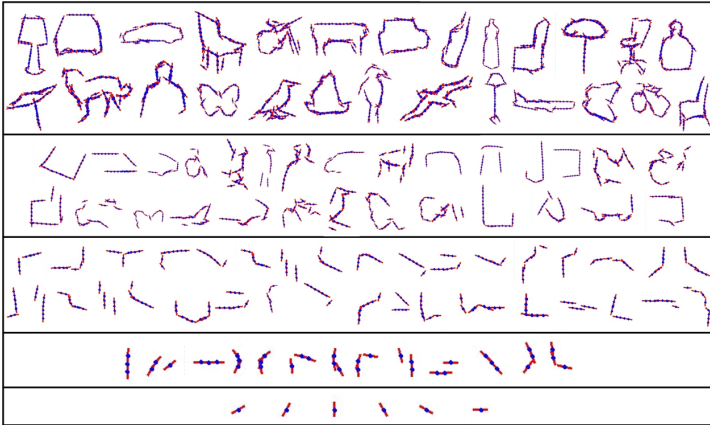


Geon theory: example of coding an object into a (partially) view-independent description

First digit- Edges: 0: Straight 1: Curved	0XXX	1XXX		
Second digit- Symmetry: 0: None 1: Reflection 2: Rotation 3: Both	X0XX	X1XX	X2XX	X3XX
Third digit- Sweep: 0: Constant Size 1: Expanding 2: Contracting 3: Both	XX0X	XX1X	XX2X	XX3X
Fourth digit- Axis: 0: Straight 1: Curved	XXX0	XXX1		
Object Relations: <0>: Smaller Than <1>: Bigger Than <2>: Above <3>: Beside <4>: Below <5>: Join End to Side <6>: Join Side to End <7>: Join both ends to side	Examples: Brick: 0300 Cylinder: 1300 Teapot: 1301<037>1310<136>1321			

Machine learning of compositional models

A 2D compositional model can also be discovered through unsupervised learning given natural image databases. A basic idea is to discover “suspicious coincidences, and then recode to remove these, in order to discover yet higher-order coincidences. Evidence for a coincidence comes from detecting when $p(A,B) \gg p(A)p(B)$.



Zhu, L., Chen, Y., & Yuille, A. (2011). Recursive Compositional Models for Vision: Description and Review of Recent Work. *Journal of Mathematical Imaging and Vision*, 41(1-2), 122–146. <http://doi.org/10.1007/s10851-011-0282-2>

Kersten, D. J., & Yuille, A. L. (2013). *Vision: Bayesian Inference and Beyond*. In J. S. Werner & L. M. Chalupa (Eds.), *The New Visual Neurosciences*. MIT Press.

Psychophysics: How sophisticated are the brain's transformation processes, between image and visual memory?

Ideal observer analysis applied to the problem of view-dependency in 3D object recognition (Liu, Knill & Kersten, 1995; Liu & Kersten, 1998)

One can imagine two quite different ways of verifying whether an unfamiliar view of an object belongs to the object or not. One way is to simply test how close the new view is to the set of stored views, without any kind of "intelligent" combination of the stored views. Given a sufficiently good representation, a simple measure of similarity could produce good recognition performance over restricted sets of views (i.e. not too much self-occlusion).

Another way is to combine the stored views in a way that reflects knowledge that they are from a 3D object, and compare the new view to the combined view. The second approach has the potential for greater generalization, and accuracy than the first.

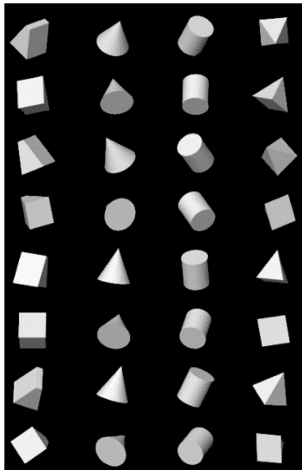
An example of the second approach would be to use the familiar views to interpolate the unfamiliar views. Given sufficient views and feature points, this latter approach has a simple mathematical realization (Ullman, 1996). An optimal verification algorithm would verify by rotating the actual 3D model of the object, projecting it to 2D and testing for an image match.

Liu et al. (1995) were able to exclude models of the first class (comparisons in 2D) and the last class (comparisons with a full 3D model) in a simple 3D classification task using ideal observer analysis. The

ideal observer technique was developed in the context of our studies of quantum efficiency in early vision.

Psychophysics: Ideal observer for the "snap shot" model of visual recognition: Discounting views

Here is an example of how to use a simple 2D image-based generative model to psychophysically address the question of what kind of image features are most effective in solving an object recognition task given varying viewpoints and visual noise.



Eight views of four objects. (See Tjan B., Braje, W., Legge, G.E. & Kersten, D. (1995) Human efficiency for recognizing 3-D objects in luminance noise. *Vision Research*, 35, 3053-3069.)

Let X = the vector describing the image data. Let O_i represent object i , where $i = 1$ to N . Suppose that O_i is represented in memory by M "snap shots" of each object, call them views (or templates) V_{ij} , where $j = 1, M$.

$$p(O_i | X) = \sum_{j=1}^M p(V_{ij} | X)$$

Set::write : Tag Times in $p(O_i | X)$ is Protected. >>

$$M p(V_{ij} | X)$$

$$= \sum_{j=1}^M \frac{p(X|V_{ij}) p(V_{ij})}{p(X)}$$

Given image data, Ideal observer chooses i that maximizes the posterior $p(O_i | X)$. If we assume that the $p(X)$ is uniform, the optimal strategy is equivalent to choosing i that maximizes:

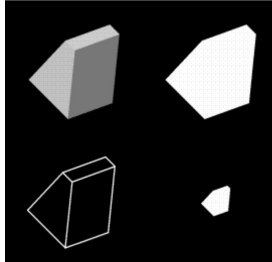
$$L(i) = \sum_{j=1}^M p(X | V_{ij}) p(V_{ij})$$

$$M p^2(X | V_{ij}) V_{ij}$$

If we assume i.i.d additive gaussian noise (as we did for the signal-known-exactly detection ideal), then

$$p(X | V_{ij}) = \frac{1}{(\sigma \sqrt{2\pi})^p} \exp\left(-\frac{1}{2\sigma^2} \|X - V_{ij}\|^2\right)$$

where the p in the exponent is the number of pixels in the image.



Tjan et al. showed that size, spatial uncertainty and detection efficiency played large roles in accounting for human object recognition efficiency. Interestingly, highest recognition efficiencies (~7.8%) were found for small silhouettes of the objects, not for line drawings. (The small silhouettes were 0.7 deg, vs. 2.4 deg for the large silhouettes).

Review of types of generative models for images

As noted above, recognition could proceed by testing whether whether models M of object fit the incoming data S , well. I.e. test $S = F(M)$? But there are many possibilities for the form of the generative model F . Below is a review of material from Lecture 7.

Generative models for images: rationale

Generative vs. discriminative models

Discriminative models for inference don't explicitly model how the image results from the object description. In Bayesian terms, an algorithm is based on:

$p(\text{object} | \text{image})$. For example, the posterior could be constructed as a look-up table: input image, check probabilities on various object descriptions, and pick the one with the biggest posterior probability.

Generative models characterize the range of variations in the image produced by an object or its representation. In Bayesian terms, algorithms explicitly model the likelihood:

$$p(\text{object} | \text{image}) \propto p(\text{image} | \text{object}) p(\text{object})$$

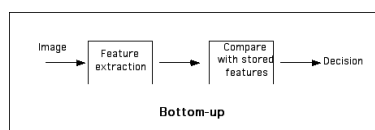
The pros for a generative are, if the visual system has built-in knowledge that can recapitulate the generative process, then recognition should be able to better generalize to novel appearances of a learned object. E.g. to deal with 3D rotations, occlusion. It can learn new objects with fewer examples.

Cons: The consensus in the field has been that modeling generative processes, especially 3D to 2D rendering, can be complex, take too much time to be practical. But see: Kulkarni, T. D., Tenenbaum, J. B., Mansinghka, V. K., & Kohli, P. (2015). Picture: A Probabilistic Programming Language for Scene Perception. Kulkarni, Tejas Dattatraya. Institute of Electrical and Electronics Engineers (IEEE)., <http://dspace.mit.edu/openaccess-disseminate/1721.1/96620>

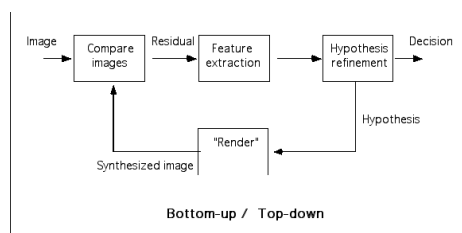
We review the types of generative models with a view to asking whether human vision can be modeled as incorporating a particular kind of generative knowledge for object recognition.

Characterize the knowledge required for inference

Feedforward procedures:



Pattern theory perspective: "analysis by synthesis"--synthesis phase explicitly incorporates generative model



Note that a top-down generative model can be used in more than one way. The above figure shows it being used to find errors in the top-down predictions. It could also be used to find consistent features.

Generative models can make it easier to conceptualize information flow: Mapping is many-to-one.

But as pointed out above, not necessarily easy to compute. For example, realistic 3D graphics rendering is computationally intensive.

Two basic concepts: Photometric & geometric variation

Two more basic concepts: 3D scene-based & 2D image-based models of geometric variation

3D Scene-based modeling: Computer graphics models

Objects & surfaces: Shape, Articulations, Material & texture. Illumination: Points and extended, Ray-tracing, Radiosity

Viewpoint/Camera

Projection geometry, homogeneous coordinates: perspective, orthographic

Application to viewpoint variation: Does human vision compensate for variations using "built-in" knowledge of 3D?

Image-based modeling

Linear intensity-based

Basis sets:

$$I = m_1 * I_1 + m_2 * I_2 + m_3 * I_3 + \dots$$

application: optics of the eye

Application in your homework: illumination variation for fixed views of an object, can be useful in object recognition

Linear geometry-based

Affine:

rigid translations, rotations, scale and shear

Application in your homework: viewpoint variation: 2D approximations to 3D variations?

Non-linear geometry-based

Morphs

Application: within-category variation for an object, or objects

finding the "average" face

Both linear and non-linear based methods raise the general question:

Does human recognition store object prototypes together with some perhaps image-based model of possible transformations to look for

a match of incoming image data with a stored template?

Modeling geometric variation: 3D scene-based modeling

Let's assume that the 2D spatial locations of certain features are stored for a given object. The math that describes how the 3D locations can be transformed and mapped on to the retinal coordinates has been known for many years, and is built into virtually all 3D graphics engines. The math will give us a handle on how to quantitatively think about matching image features to stored memory.

Let's look at the math. There are four basic types of transformations: rotation, scale, translation, and projection. First rotations. Then we'll put rotations together with the other transformations using homogeneous coordinates.

Representing Rotations

See `Rotate[]` and `RotationMatrix[]` for built-in Mathematica functions for doing rotations.

Euler angles

Euler angles are a standard way of representing rotations of a rigid body.

A rotation specified by the Euler angles ψ , θ , and ϕ can be decomposed into a sequence of three successive rotations: first by angle ψ about the z axis, the second by angle θ about the x axis, and the third about the z axis (again) by angle ϕ . The angle θ is restricted to the range 0 to π .

$$\text{In[2]:= RotationMatrix3D}[\psi, \theta, \phi] := \begin{pmatrix} \text{Cos}[\phi] \text{Cos}[\psi] - \text{Cos}[\theta] \text{Sin}[\phi] \text{Sin}[\psi] & \text{Cos}[\theta] \text{Cos}[\psi] \text{Sin}[\phi] & \text{Sin}[\theta] \text{Sin}[\phi] \\ -\text{Cos}[\psi] \text{Sin}[\phi] - \text{Cos}[\theta] \text{Cos}[\phi] \text{Sin}[\psi] & \text{Cos}[\theta] \text{Cos}[\phi] \text{Cos}[\psi] - \text{Sin}[\theta] \text{Sin}[\psi] & \text{Cos}[\phi] \text{Sin}[\theta] \\ \text{Sin}[\theta] \text{Sin}[\psi] & -\text{Cos}[\psi] \text{Sin}[\theta] & \text{Cos}[\theta] \end{pmatrix}$$

In[3]:= `RotationMatrix3D[ψ, θ, φ] // MatrixForm`

Out[3]//MatrixForm=

$$\begin{pmatrix} \text{Cos}[\phi] \text{Cos}[\psi] - \text{Cos}[\theta] \text{Sin}[\phi] \text{Sin}[\psi] & \text{Cos}[\theta] \text{Cos}[\psi] \text{Sin}[\phi] + \text{Cos}[\phi] \text{Sin}[\psi] & \text{Sin}[\theta] \text{Sin}[\phi] \\ -\text{Cos}[\psi] \text{Sin}[\phi] - \text{Cos}[\theta] \text{Cos}[\phi] \text{Sin}[\psi] & \text{Cos}[\theta] \text{Cos}[\phi] \text{Cos}[\psi] - \text{Sin}[\theta] \text{Sin}[\psi] & \text{Cos}[\phi] \text{Sin}[\theta] \\ \text{Sin}[\theta] \text{Sin}[\psi] & -\text{Cos}[\psi] \text{Sin}[\theta] & \text{Cos}[\theta] \end{pmatrix}$$

In[4]:= `RotationMatrix3D[ψ, θ, 0] // MatrixForm`

Out[4]//MatrixForm=

$$\begin{pmatrix} \text{Cos}[\psi] & \text{Sin}[\psi] & 0 \\ -\text{Sin}[\psi] & \text{Cos}[\psi] & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

In[5]:= `RotationMatrix3D[0, θ, 0] // MatrixForm`

Out[5]//MatrixForm=

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & \text{Cos}[\theta] & \text{Sin}[\theta] \\ 0 & -\text{Sin}[\theta] & \text{Cos}[\theta] \end{pmatrix}$$

In[6]:= `RotationMatrix3D[0, 0, φ] // MatrixForm`

Out[6]//MatrixForm=

$$\begin{pmatrix} \text{Cos}[\phi] & \text{Sin}[\phi] & 0 \\ -\text{Sin}[\phi] & \text{Cos}[\phi] & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

Homogeneous coordinates

Rotation and scaling can be done by linear matrix operations in three-space. Translation and perspective transformations do not have a three dimensional matrix representation. By going from three dimensions to four dimensional coordinates, all four of the above basic operations can be represented within the formalism of matrix multiplication.

Homogeneous coordinates are defined by: $\{xw, yw, zw, w\}$, (w not equal to 0). To get from homogeneous coordinates to three-space coordinates, $\{x,y,z\}$, divide the first three homogeneous coordinates by the fourth, $\{w\}$. For more information, see references on 3D graphics, e.g. Foley, J., van Dam, A., Feiner, S., & Hughes, J. (1990).

The rotation and translation matrices can be used to describe object or eye-point changes of position. The scaling matrix allows you to squash or expand objects in any of the three directions. Any combina-

tion of the matrices can be multiplied together or concatenated. But remember, matrices do not in general commute, so the order is important. The translation, rotation, and perspective transformation matrices can be concatenated to describe general 3-D to 2-D perspective mappings.

We will use these definitions later when we develop a theory for computing structure from motion, spatial layout and direction of heading.

```
In[7]:= XRotationMatrix[theta_] := {{1, 0, 0, 0},
    {0, Cos[theta], Sin[theta], 0}, {0, -Sin[theta], Cos[theta], 0}, {0, 0, 0, 1}};
YRotationMatrix[theta_] := {{Cos[theta], 0, -Sin[theta], 0},
    {0, 1, 0, 0}, {Sin[theta], 0, Cos[theta], 0}, {0, 0, 0, 1}};
ZRotationMatrix[theta_] := {{Cos[theta], Sin[theta], 0, 0},
    {-Sin[theta], Cos[theta], 0, 0}, {0, 0, 1, 0}, {0, 0, 0, 1}};
ScaleMatrix[sx_, sy_, sz_] := {{sx, 0, 0, 0}, {0, sy, 0, 0}, {0, 0, sz, 0}, {0, 0, 0, 1}};
TranslateMatrix[x_, y_, z_] := {{1, 0, 0, 0}, {0, 1, 0, 0}, {0, 0, 1, 0}, {x, y, z, 1}};
ThreeDToHomogeneous[vec_] := Append[vec, 1];
HomogeneousToThreeD[vec_] := Drop[ $\frac{\text{vec}}{\text{vec}[[4]}}$ , -1];
ZProjectMatrix[focal_] :=
    {{1, 0, 0, 0}, {0, 1, 0, 0}, {0, 0, 0, -N[ $\frac{1}{\text{focal}}$ ]}, {0, 0, 0, 1}};
ZOrthographic[vec_] := Take[vec, 2];
```

Translation by $\{d_x, d_y, d_z\}$ can be found by applying the matrix

```
In[16]:= Clear[d];
TranslateMatrix[dx, dy, dz] // MatrixForm
```

```
Out[17]//MatrixForm=

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ d_x & d_y & d_z & 1 \end{pmatrix}$$

```

```
In[18]:= {x, y, z, 1}.TranslateMatrix[dx, dy, dz]
```

```
Out[18]= {x + dx, y + dy, z + dz, 1}
```

```
to {x,y,z,1}
```

$$(x + d_x, y + d_y, z + d_z, 1) = (x, y, z, 1) \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ d_x & d_y & d_z & 1 \end{pmatrix}$$

The scaling matrix is:

```
In[19]:= ScaleMatrix[sx, sy, sz] // MatrixForm
```

```
Out[19]//MatrixForm=
```

$$\begin{pmatrix} s_x & 0 & 0 & 0 \\ 0 & s_y & 0 & 0 \\ 0 & 0 & s_z & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

There are three matrices for general rotation:

z-axis (moving the positive x-axis towards the positive y-axis)

```
In[20]:= ZRotationMatrix[θ] // MatrixForm
```

```
Out[20]//MatrixForm=
```

$$\begin{pmatrix} \cos[\theta] & \sin[\theta] & 0 & 0 \\ -\sin[\theta] & \cos[\theta] & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

x-axis (moving the positive y towards the positive z)

```
In[21]:= XRotationMatrix[θ] // MatrixForm
```

```
Out[21]//MatrixForm=
```

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos[\theta] & \sin[\theta] & 0 \\ 0 & -\sin[\theta] & \cos[\theta] & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

y-axis (moving positive z towards positive x):

```
In[22]:= YRotationMatrix[θ] // MatrixForm
```

```
Out[22]//MatrixForm=
```

$$\begin{pmatrix} \cos[\theta] & 0 & -\sin[\theta] & 0 \\ 0 & 1 & 0 & 0 \\ \sin[\theta] & 0 & \cos[\theta] & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

Perspective

Perspective transformation is the only one that requires extracting the three-space coordinates by dividing the homogeneous coordinates by the fourth component w . The projection plane is the x - y plane, and the focal point is at $z = d$. Then $\{x, y, z, 1\}$ maps onto $\{x, y, 0, -z/d + 1\}$ by the following transformation:

```
In[23]:= Clear[d]
ZProjectMatrix[d] // MatrixForm
```

```
Out[24]/MatrixForm=

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & -\frac{1}{d} \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

```

After normalization, the image coordinates $\{x',y',z'\}$ are read from:

$$(x', y', z', 1) = \left(\frac{xd}{d-z}, \frac{yd}{d-z}, 0, 1 \right)$$

The steps can be seen here:

```
In[25]:= Clear[x, y, z, d]
{x, y, z, 1}.ZProjectMatrix[d]
{x, y, z, 1}.ZProjectMatrix[d] / %[[4]]
HomogeneousToThreeD[{x, y, z, 1}.ZProjectMatrix[d]]
Simplify[ZOrthographic[HomogeneousToThreeD[{x, y, z, 1}.ZProjectMatrix[d]]]]
```

```
Out[26]= {x, y, 0, 1 - \frac{z}{d}}
```

```
Out[27]= \left\{ \frac{x}{1 - \frac{z}{d}}, \frac{y}{1 - \frac{z}{d}}, 0, 1 \right\}
```

```
Out[28]= \left\{ \frac{x}{1 - \frac{z}{d}}, \frac{y}{1 - \frac{z}{d}}, 0 \right\}
```

```
Out[29]= \left\{ \frac{d x}{d - z}, \frac{d y}{d - z} \right\}
```

The matrix for orthographic projection has $d \rightarrow$ infinity.

```
In[30]:= Limit[ZOrthographic[HomogeneousToThreeD[{x, y, z, 1}.ZProjectMatrix[d]]], d -> \infty]
```

```
Out[30]= {x, y}
```

The perspective transformation is the only singular matrix in the above group. This means that, unlike the others its operation is not invertible. Given the image coordinates, the original scene points cannot be determined uniquely.

Example: transforming, projecting a 3D object

We are going to generate a "view" of random 3D object. Imagine that you've seen this view (threeDtemplate) and stored it in memory. Later you get a view (newvertices) of either the same object or a different one, and you want to check if it is the same object as before. You need to make some kind of comparison test.

We'll keep it simple and do orthographic projection.

```
In[31]:= orthoproject[x_] := Delete[x, Table[{i, 3}, {i, 1, Length[x]}]];
```

Define 3D target object - Wire with randomly positioned vertices

```
In[32]:= threeDtemplate = Table[{RandomReal[], RandomReal[], RandomReal[]}, {5}];
```

First view

View from along Z-direction

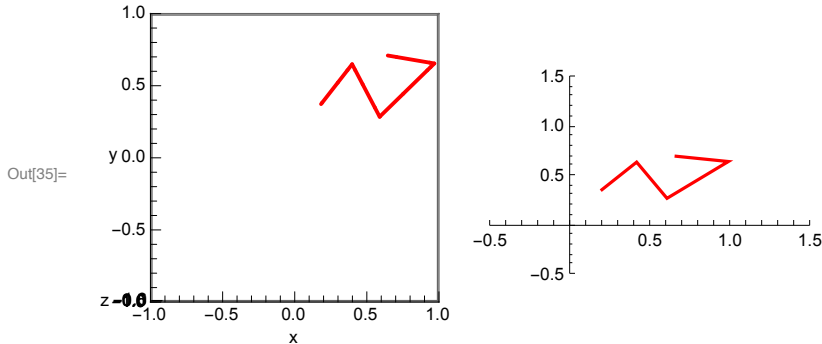
```
In[33]:= lines = Partition[threeDtemplate, 2, 1];
fv3d = Graphics3D[{Thick, Red, Line[lines]}, ViewPoint -> {0, 0, 100},
  PlotRange -> {{-1, 1}, {-1, 1}, {-1, 1}}, AspectRatio -> 1, Axes -> True,
  AxesLabel -> {"x", "y", "z"}, ImageSize -> Small, PreserveImageOptions -> True];
```

ListPlot view

We can also do the projection ourselves:

```
In[34]:= ovg = ListPlot[orthoproject[threeDtemplate], Joined -> True, PlotStyle ->
  {Thickness[0.01], RGBColor[1, 0, 0]}, PlotRange -> {{-.5, 1.5}, {-.5, 1.5}}];
```

```
In[35]:= GraphicsRow[{fv3d, ovg}]
```



New View

We pick an arbitrary view of the above object.

Use Homogeneous coordinates

```
In[36]:= swidth = 1.0; sheight = 1.0; slength = 1.0; d = 0;
```

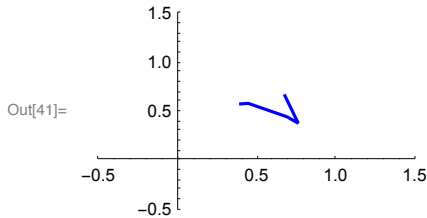
```
In[37]:= homovertrices = Transpose[Map[ThreeDToHomogeneous, threeDtemplate]];
newtransformMatrix = TranslateMatrix[.3, 0, 0].XRotationMatrix[N[ $\frac{\pi}{2} * .3$ ]].
  YRotationMatrix[N[ $-\frac{\pi}{2} * .2$ ]].ScaleMatrix[swidth, sheight, slength];
```

```
In[39]:= temp = N[newtransformMatrix.homovertrices];
```

Take a look at the new view

```
In[40]:= newvertices = Map[HomogeneousToThreeD, Transpose[temp]];
```

```
In[41]:= ListPlot[orthoproject[newvertices], Joined -> True,
  PlotStyle -> {Thickness[0.01], RGBColor[0, 0, 1]},
  PlotRange -> {{-.5, 1.5}, {-.5, 1.5}}, ImageSize -> Small]
```



- ▶ 1. Exercise: look at new view by coding the orthographic projection yourself

Modeling geometric variation: linear 2D image, geometry-based modeling as an approximation to 3D scene-based variation

Suppose that we encounter a new view of the 3D object, i.e. from some new arbitrary viewpoint. This new viewpoint can be modeled as a 3D rotation and translation of the object.

If we want to see if the new and old images are of the same object, we could try to rotate a 3D representation of the object. But this would require knowledge of 3D.

Alternatively, if one projects a rotation in 3D onto a 2D view, we can try to approximate the rotation by a 2D affine transformation. A 2D affine transformation is a simple 2D operation, perhaps it is sufficient to account for the generalization of familiar to unfamiliar views?

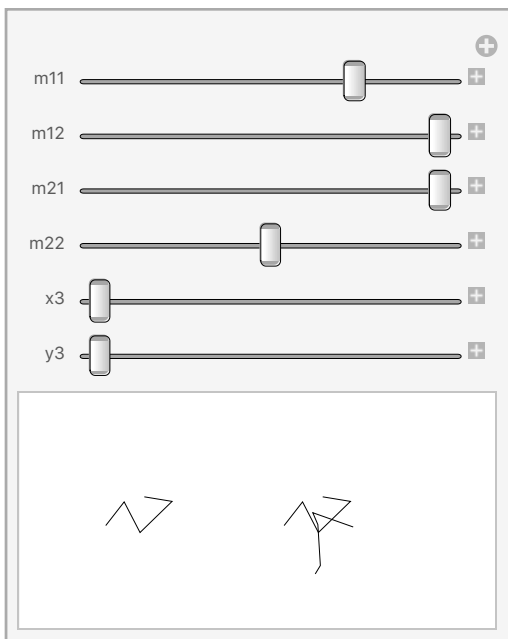
Affine transformation preserves parallel lines.

We know that rotations, scale and shear transformations will preserve parallel lines. So will translations. It is not immediately apparent, that any matrix operation is an affine transformation, although one has to remember that translations are not represented by matrix operations unless one goes to homogeneous coordinates. Lecture 7 had a simple demo of the parallel line preservation for transformations of a square.

Try to find values of a 2D matrix ($M2=\{\{m11,m12\},\{m21,m22\}\}$) and 2D translations ($x3,y3$) that bring newvertices as close as possible to the threeDtemplate stored in memory.

```
In[42]:= Manipulate[
  x1 = orthoproject[threeDtemplate];
  M2 = {{m11, m12}, {m21, m22}};
  x2 = (M2.#1 &) /@ orthoproject[newvertices];
  x2b = # + {x3, y3} & /@ x2;
  GraphicsRow[{Graphics[Line[x1], PlotRange → {{-.5, 1.5}, {-.5, 1.5}}],
    Graphics[{Line[x1], Line[x2b]}, PlotRange → {{-.5, 1.5}, {-.5, 1.5}}]},
  ImageSize → Small],
  {{m11, 1}, -2, 2}, {{m12, 2}, -2, 2}, {{m21, 2}, -2, 2},
  {{m22, 0}, -2, 2}, {x3, -1, 1}, {y3, -1, 1}]
```

Out[42]=



Compute closest least squares affine match with translation

```
In[43]:= aff = {{aa, bb}, {cc, dd}}; tra = {ff, gg};
errorsum :=
  Apply[Plus, Flatten[ (# + tra & /@ (aff.#1 &) /@ orthoproject[newvertices] -
    orthoproject[threeDtemplate]) ^2]];
temp = FindMinimum[errorsum, {aa, .8}, {bb, .2}, {cc, .16}, {dd, .8},
  {ff, 0.0}, {gg, 0.0}, MaxIterations -> 200];
minvals = Take[temp, -1][[1]]; minerr = Take[temp, 1][[1]];
naff = aff /. minvals; ntra = tra /. minvals;
minerr
```

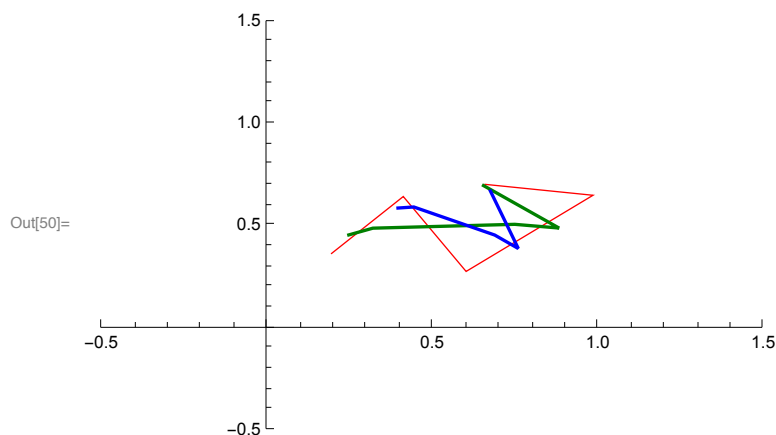
```
Out[48]= 0.154175
```

Check match with estimated view

```
In[49]:= estim = naff.Transpose[orthoproject[newvertices]] + ntra;
```

Plot first original view, new view and the affine estimate of the first from the new

```
In[50]:= evg = ListPlot[
  {orthoproject[threeDtemplate], Transpose[estim], orthoproject[newvertices]},
  Joined -> True, PlotStyle -> {{Thickness[0.002], RGBColor[1, 0, 0]},
    {Thickness[0.005], RGBColor[0, .5, 0]},
    {Thickness[0.005], RGBColor[0, 0, 1]}},
  PlotRange -> {{-.5, 1.5}, {-.5, 1.5}}
```



Liu & Kersten (1998) compared human recognition performance with 2D affine observers. The targets were paper-clip like objects as above, except thicker with some shading. Human performance was somewhat better than the affine observer, suggesting that people can incorporate additional 3D information, perhaps from the shading/occlusion information, together with a "smarter" model.

Appendix: Neuropsychological and neurophysiological studies

Neuropsychological Studies

Category-specific breakdowns

Inferomedial occipito-temporal region, (right hemi), fusiform and lingual gyri--> prosopagnosia. Can recognize other objects (even with comparable structural complexity), and can recognize a face as a face, and can name its parts.

...but is it a problem with individuation in a class? Evidence suggesting prosopagnosics have a problem distinguishing fruits, playing cards, autos, etc.. Bird-watcher lost ability. Farmer couldn't identify his cows.

Damasio's patients could recognize horses, owls, elephants, but had problems with dollar sign, British pound sign, musical clef. --> perhaps a problem with inter-category discriminations (subordinate-level), rather than complexity per se.

Corroboration--patient with car agnosia could still identify ambulance and fire engine (distinct entry point attributes)

BUT, prosopagnosia does seem sometimes to occur without any of the subordinate-level deficit. Patients impaired for living, but not non-living things.

<<20 questions and recognition>>

Summary: Two types of visual memory:

recognition that involves representing and distinguishing prototypes

<<Different prototypes in different IT hypercolumns?>>

recognition that involves distinguishing deviations between members with the same prototype (inferomedial occipito-temporal)

<<processing within hypercolumn?>>

Deficits in recognizing facial expressions

Dissociation between face recognition and recognizing facial expressions.

Some prosopagnosics can't recognize an individual face, but can recognize the expression.

Damasio reports bilateral amygdala lesion patient could recognize individual faces, but did not do well with expressions of happiness, surprise, fear, anger, etc.. Monkeys too (Weiskrantz, 1956)

Metamorphopsia with faces. Another patient experiences metamorphopsia with objects other than faces.

Visuomotor

DF

Electrophysiological Studies

V1-> V2 -> V4 -> IT-> TEO (PIT) -> TE
not strictly serial

V2, V3, V4, corpus callosum -> IT

TE, TEO connected to thalamus, hypothalamus,...

Object information might even skip IT and go to limbic structures or striatum...

>abstract categorizations (with high cue validity) perhaps possible even with damage to TE

Physiological properties of IT neurons

Physiological properties of IT neurons

Gross. IT as last exclusive visual area.

Posterior TEO, cells similar to V4, visuotopic, repres. contralateral vis. field, rfs larger than V4. (small as 1.5 - 2.5 deg)

anterior TE, complex stimuli required. TE not visuotopic, large ipsi, contra or bilat. rfs.

30 to 50 deg rfs.

Cells often respond more vigorously to Fovea stimulation

Shape selectivity (some in V4), lots in IT. natural objects, Walsh functions, faces, hands.

Invariance? Rare to find size or position constancy--but selectivity falls off slowly over size and position.

Thus in this sense roughly 50% of cells show size and position invariance.

Cue invariant--motion, texture or luminance defined shape boundaries. BUT, contrast polarity sensitive. >>shape from shading?

Two mechanisms? 1) prototypes of objects that can be decomposed into parts.

parts important.

2) holistic, configurational. Part features not useful for discrimination, but whole is.

Combination encoding

Tanaka & modules for similar shapes, columnar organization. \

>1300 prototype modules?? RBC?

Sufficient for representing an exemplar of a category? Or when holistic information is required?

L&S suggest combination encoding not used for holistic representation. Evidence: Many cells in TE and STS code overall shape of biologically important objects--not features or parts. Novel wire objects too.

Selectivity for biologically important stimuli

Face cells - TEa, TEm, STS, amygdala, inf. convexity of prefrontal cortex.

Some cells like features (e.g. eyes). Other like the whole face, or face-view, or even highly selective for face-gaze angle, head direction, and body posture.

Face cells, invariant over size and position, less so over orientation--upright preferred.

Face identity cells in IT,

but facial expression, gaze direction, and vantage point in STS

PET, posterior fusiform gyrus for face matching, gender disc.

mid-fusiform for unique face

IT cells for whole human body, mostly viewer centered cells. 20% holistic

Configurational selectivity for novel objects

L et al., and L&S's work. on wires, etc.
 ant. medial temporal sulcus
 view-selective "blurred templates"
 enantiomorphic views undistinguished
 many showed broad size tuning
 Action-related
 MT -> parietal MST, FST, LIP, 7,
 LIP cells sensitive to grasp shape of hand

Compute closest least squares affine match without translation

```
In[51]:= naff2 = Transpose[orthoproject[threeDtemplate]].  
          PseudoInverse[Transpose[orthoproject[newvertices]]]
```

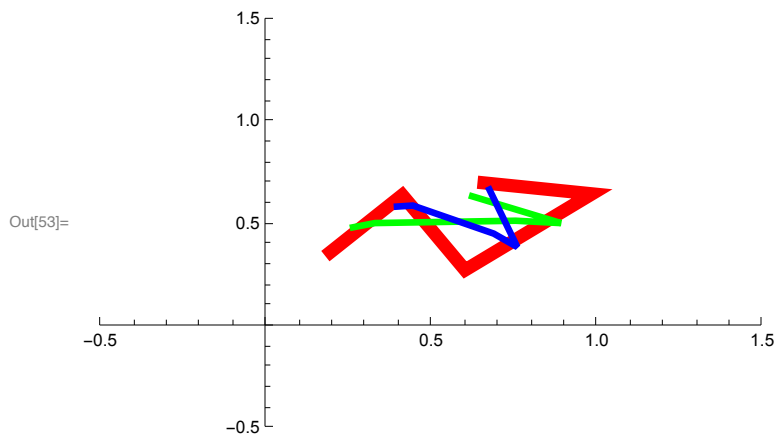
```
Out[51]:= {{1.44339, -0.523846}, {0.371013, 0.578403}}
```

Check match with estimated view

```
In[52]:= estim2 = naff2.Transpose[orthoproject[newvertices]];
```

Plot familiar view, new view and the affine estimate of the old from the new

```
In[53]:= evg = ListPlot[  
  {orthoproject[threeDtemplate], Transpose[estim2], orthoproject[newvertices]},  
  PlotJoined -> True, PlotStyle -> {{Thickness[0.02], RGBColor[1, 0, 0]},  
  {Thickness[0.01], RGBColor[0, 1, 0]},  
  {Thickness[0.01], RGBColor[0, 0, 1]}},  
  PlotRange -> {{-0.5, 1.5}, {-0.5, 1.5}}
```



Test set of newvertices and threeDtemplate

$$\text{In[54]:= } \left(\begin{array}{ccc} \text{*threeDtemplate=} & \begin{pmatrix} 0.23981762582649485 & 0.14312418380466885 & 0.03003120544761813 \\ 0.2624091279705781 & 0.4565009537332048 & 0.1221875974954246 \\ 0.019392922865028396 & 0.016530310373452352 & 0.5906147114395374 \\ 0.06481020981636326 & 0.6548152420848915 & 0.40459291550719 \\ 0.6422482206653176 & 0.7719461816974882 & 0.22053936016974654 \end{pmatrix} & \end{array} \right) \text{*)}$$

$$\text{In[55]:= } \left(\begin{array}{ccc} \text{*newvertices=} & \begin{pmatrix} 0.2215818503538964 & 0.09974436717830631 & -0.09854211812818665 \\ 0.26452283137288907 & 0.3891455135818127 & -0.16199300398045482 \\ 0.18952784909991596 & 0.25183584120022917 & 0.45991480775746235 \\ 0.17676551774440416 & 0.7093221832702079 & 0.026256226849368618 \\ 0.5640697249874365 & 0.5756717475918378 & -0.28280208011255054 \end{pmatrix} & \end{array} \right) \text{*)}$$

References

Recognition

- Biederman, I. (1987). Recognition-by-components: A theory of human image understanding. *Psychological Review*, 94, 115-147.
- Bülthoff, H. H., & Edelman, S. (1992). Psychophysical support for a two-dimensional view interpolation theory of object recognition. *Proc. Natl. Acad. Sci. USA*, 89, 60-64.
- Clark, J. J., & Yuille, A. L. (1990). *Data Fusion for Sensory Information Processing*. Boston: Kluwer Academic Publishers.
- David, C., & Zucker, S. W. (1989). Potentials, Valleys, and Dynamic Global Coverings (TR-CIM 98-1): McGill Research Centre for Intelligent Machines, McGill University.
- DiCarlo, J. J., Zoccolan, D., & Rust, N. C. (2012). How does the brain solve visual object recognition? *Neuron*, 73(3), 415-434. <http://doi.org/10.1016/j.neuron.2012.01.010>
- Field, D. J., Hayes, A., & Hess, R. F. (1993). Contour integration by the human visual system: evidence for a local "association field". *Vision Research*, 33, 173-193.
- Kersten, D. J. (1991). Transparency and the Cooperative Computation of Scene Attributes. In M. Landy & A. Movshon (Eds.), *Computational Models of Visual Processing*, (pp. 209-228). Cambridge, Massachusetts: M.I.T. Press.
- Kersten, D. & Madarasmı, S. (1995). The Visual Perception of Surfaces, their Properties, and Relationships. In I. J. Cox, P. Hansen, & B. Julesz (Ed.), *Partitioning Data Sets: With applications to psychology, vision and target tracking*. (pp. 373-389). American Mathematical Society.
- Kersten, D. (1999). High-level vision as statistical inference. In M. S. Gazzaniga (Ed.), *The New Cognitive Neurosciences -- 2nd Edition* (pp. 353-363). Cambridge, MA: MIT Press.
- Kersten, D., & Schrater, P. W. (2000). Pattern Inference Theory: A Probabilistic Approach to Vision. In R. Mausfeld & D. Heyer (Eds.), *Perception and the Physical World*. Chichester: John Wiley & Sons, Ltd.
- Liu, Z., Knill, D. C. & Kersten, D. (1995). Object Classification for Human and Ideal Observers. *Vision Research*, 35, 549-568.

- Liu, Z., & Kersten, D. (1998). 2D observers for 3D object recognition? In *Advances in Neural Information Processing Systems* Cambridge, Massachusetts: MIT Press.
- Logothetis, N. K., Pauls, J., Bulthoff, H. H. & Poggio, T. (1994). View-dependent object recognition by monkeys. *Current Biology*, 4 No 5, 401-414.
- Logothetis, N. K., & Sheinberg, D. L. (1996). Visual Object Recognition. *Annual Review of Neuroscience*, 19, 577-621.
- Mohan, R. (1989). *Perceptual organization for computer vision (IRIS 254)*: University of Southern California.
- Poggio, T. & Edelman, S. (1990). A network that learns to recognize three-dimensional objects. *Nature*, 343, 263-266.
- Riesenhuber, M., & Poggio, T. (1999). Hierarchical models of object recognition in cortex. *Nature Neuroscience*, 2, 1019–1025.
- Rock, I. & Di Vita, J. (1987). A case of viewer-centered object perception. *Cognitive Psychology*, 19, 280-293.
- Sharon, E., Galun, M., Sharon, D., Basri, R., & Brandt, A. (2006). Hierarchy and adaptivity in segmenting visual scenes. *Nature*, 442(7104), 810–813. doi:10.1038/nature04977
- Shashua, A., & Ullman, S. (1988,). Structural Saliency: The detection of globally salient structures using a locally connected network. Paper presented at the 2nd International Conference on Computer Vision, Washington, D.C.
- Shi, J., & Malik, J. (2000). Normalized cuts and image segmentation. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 22(8), 888–905.
- Tanaka, K. (1996). Inferotemporal cortex and object vision. *Annual Review of Neuroscience*, 19, 109-139.
- Tarr, M. J., & Bülthoff, H. H. (1995). Is human object recognition better described by geon-structural-descriptions or by multiple-views? *Journal of Experimental Psychology: Human Perception and Performance*, 21(6), 1494-1505.
- Troje, N. F., & Kersten, D. (1999). Viewpoint dependent recognition of familiar faces. *Perception*, 28(4), 483 - 487.
- Ullman, S. (1996). *High-level Vision: Object Recognition and Visual Cognition*. Cambridge, Massachusetts: MIT Press.
- Ullman, S., Vidal-Naquet, M., & Sali, E. (2002). Visual features of intermediate complexity and their use in classification. *Nat Neurosci*, 5(7), 682-687.
- Vetter, T., Poggio, T., & Bülthoff, H. H. (1994). The importance of symmetry and virtual views in three-dimensional object recognition. *Current Biology*, 4(1), 18-23.
- Yamins, D. L. K., Hong, H., Cadieu, C. F., Solomon, E. A., Seibert, D., & DiCarlo, J. J. (2014). Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the National Academy of Sciences of the United States of America*, 111(23), 8619–8624.
- Yuille, A. (2011). Towards a theory of compositional learning and encoding of objects. *Computer Vision Workshops (ICCV Workshops), 2011 IEEE International Conference on*, 1448–1455.