

**Recognition wrap-up  
&  
Self-supervised representation  
learning**

Kristen Grauman  
 UT-Austin  
 Wed Sept 20, 2017

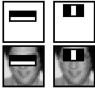

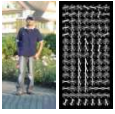
Announcements

- Assignment 1 due Sept 22 11:59 pm on Canvas
- Hw2 is out and due Wed Oct 11
- Next week: CNN hands-on tutorial with Ruohan Gao and Tushar Nagarajan
  - Bring laptop
  - Set up your TACC portal account in advance

Outline

- Last time
  - Spatial verification for instance recognition
  - Recognizing categories
- Today
  - Wrap up on categories/classifiers
  - Self-supervised learning
  - External papers & assigned paper discussion
    - Shuffle and Learn (Yu-Chuan)
    - Colorization (Keivaun)
    - Curious Robot (Ginevra)
  - Experiment
    - Network dissection (Thomas and Wonjoon)

Last time: Three landmark case studies for image classification

		
Boosting + face detection	NN + scene Gist classification	SVM + person detection
Viola & Jones	e.g., Hays & Efros	e.g., Dalal & Triggs

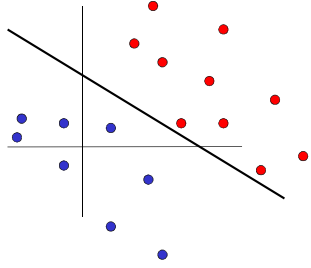
Slide credit: Kristen Grauman

Last time

- Intro to categorization problem
- Object categorization as discriminative classification
  - Boosting + fast face detection example
  - Nearest neighbors + scene recognition example
  - Support vector machines + pedestrian detection example
    - Pyramid match kernels, spatial pyramid match
  - Convolutional neural networks + ImageNet example

Linear classifiers

---



### Linear classifiers

- Find linear function to separate positive and negative examples

$x_i$  positive:  $x_i \cdot w + b \geq 0$   
 $x_i$  negative:  $x_i \cdot w + b < 0$

Which line is best?

### Support Vector Machines (SVMs)

- Discriminative classifier based on *optimal separating hyperplane*
- Maximize the *margin* between the positive and negative training examples

### Support vector machines

- Want line that maximizes the margin.

$x_i$  positive ( $y_i = 1$ ):  $x_i \cdot w + b \geq 1$   
 $x_i$  negative ( $y_i = -1$ ):  $x_i \cdot w + b \leq -1$

For support, vectors,  $x_i \cdot w + b = \pm 1$

Support vectors      Margin

C. Burges, [A Tutorial on Support Vector Machines for Pattern Recognition](#), Data Mining and Knowledge Discovery, 1998

### Support vector machines

- Want line that maximizes the margin.

$x_i$  positive ( $y_i = 1$ ):  $x_i \cdot w + b \geq 1$   
 $x_i$  negative ( $y_i = -1$ ):  $x_i \cdot w + b \leq -1$

For support, vectors,  $x_i \cdot w + b = \pm 1$

Distance between point and line:  $\frac{|x_i \cdot w + b|}{\|w\|}$

For support vectors:

$$\frac{w^T x + b}{\|w\|} = \frac{\pm 1}{\|w\|} \quad M = \left| \frac{1}{\|w\|} - \frac{-1}{\|w\|} \right| = \frac{2}{\|w\|}$$

Support vectors      Margin M

### Support vector machines

- Want line that maximizes the margin.

$x_i$  positive ( $y_i = 1$ ):  $x_i \cdot w + b \geq 1$   
 $x_i$  negative ( $y_i = -1$ ):  $x_i \cdot w + b \leq -1$

For support, vectors,  $x_i \cdot w + b = \pm 1$

Distance between point and line:  $\frac{|x_i \cdot w + b|}{\|w\|}$

Therefore, the margin is  $2 / \|w\|$

Support vectors      Margin M

### Finding the maximum margin line

- Maximize margin  $2/\|w\|$
- Correctly classify all training data points:
  - $x_i$  positive ( $y_i = 1$ ):  $x_i \cdot w + b \geq 1$
  - $x_i$  negative ( $y_i = -1$ ):  $x_i \cdot w + b \leq -1$

*Quadratic optimization problem:*

Minimize  $\frac{1}{2} w^T w$

Subject to  $y_i(w \cdot x_i + b) \geq 1$

### Finding the maximum margin line

- Solution:  $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$

learned weight

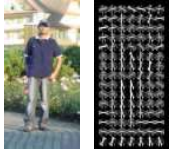
Support vector

### Finding the maximum margin line

- Solution:  $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$   
 $b = y_i - \mathbf{w} \cdot \mathbf{x}_i$  (for any support vector)
- Classification function:  
 $f(x) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$   
 $= \text{sign}\left(\sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b\right)$

C. Burges, A Tutorial on Support Vector Machines for Pattern Recognition, Data Mining and Knowledge Discovery

### Person detection with HoG's & linear SVM's

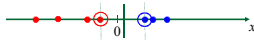
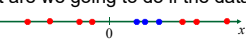
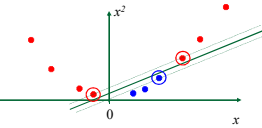


- Map each grid cell in the input window to a histogram counting the gradients per orientation.
- Train a linear SVM using training set of pedestrian vs. non-pedestrian windows.

Dalal & Triggs, CVPR 2005

Code available: <http://pascal.inrialpes.fr/soft/oltp/>

### Non-linear SVMs

- Datasets that are linearly separable with some noise work out great: 
- But what are we going to do if the dataset is just too hard? 
- How about... mapping data to a higher-dimensional space: 

### Nonlinear SVMs

- The kernel trick*: instead of explicitly computing the lifting transformation  $\phi(x)$ , define a kernel function  $K$  such that
 
$$K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$$
- This gives a nonlinear decision boundary in the original feature space:
 
$$\sum_i \alpha_i y_i K(x_i, \mathbf{x}) + b$$

### Example

2-dimensional vectors  $\mathbf{x} = [x_1 \ x_2]$ ;  
 let  $K(x_i, x_j) = (1 + x_i^T x_j)^2$

Need to show that  $K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$ :

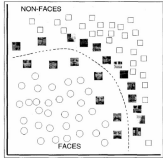
$$\begin{aligned}
 K(x_i, x_j) &= (1 + x_i^T x_j)^2 \\
 &= 1 + x_{i1}^2 x_{j1}^2 + 2 x_{i1} x_{j1} x_{i2} x_{j2} + x_{i2}^2 x_{j2}^2 + 2 x_{i1} x_{j1} + 2 x_{i2} x_{j2} \\
 &= [1 \ x_{i1}^2 \ \sqrt{2} x_{i1} x_{i2} \ x_{i2}^2 \ \sqrt{2} x_{i1} \ \sqrt{2} x_{i2}]^T \\
 &\quad [1 \ x_{j1}^2 \ \sqrt{2} x_{j1} x_{j2} \ x_{j2}^2 \ \sqrt{2} x_{j1} \ \sqrt{2} x_{j2}] \\
 &= \phi(x_i)^T \phi(x_j), \\
 &\text{where } \phi(x) = [1 \ x_1^2 \ \sqrt{2} x_1 x_2 \ x_2^2 \ \sqrt{2} x_1 \ \sqrt{2} x_2]
 \end{aligned}$$

### Examples of kernel functions

- Linear:  $K(x_i, x_j) = x_i^T x_j$
- Gaussian RBF:  $K(x_i, x_j) = \exp(-\frac{\|x_i - x_j\|^2}{2\sigma^2})$
- Histogram intersection:  $K(x_i, x_j) = \sum_k \min(x_i(k), x_j(k))$

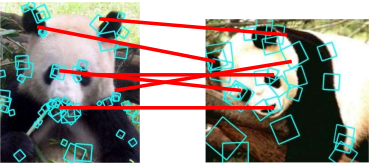
### SVMs for recognition

1. Define your representation for each example.
2. Select a kernel function.
3. Compute pairwise kernel values between labeled examples
4. Use this "kernel matrix" to solve for SVM support vectors & weights.
5. To classify a new example: compute kernel values between new input and support vectors, apply weights, check sign of output.



Kristen Grauman

### What about a matching kernel?




$X = \{\vec{x}_1, \dots, \vec{x}_m\}$      $Y = \{\vec{y}_1, \dots, \vec{y}_n\}$

Local feature correspondence useful similarity measure for generic object categories

Kristen Grauman

### Partially matching sets of features



Optimal match:  $O(m^3)$   
 Greedy match:  $O(m^2 \log m)$   
 Pyramid match:  $O(m)$

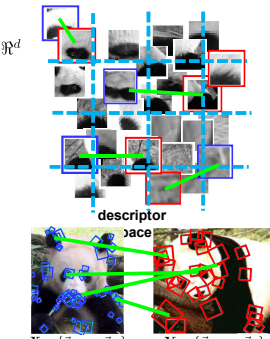
$X = \{\vec{x}_1, \dots, \vec{x}_m\}$      $Y = \{\vec{y}_1, \dots, \vec{y}_n\}$     (m=num pts)

$\min_{\pi: X \rightarrow Y} \sum_{x_i \in X} \|x_i - \pi(x_i)\|$     hate matching kernel that makes it practical to compare large sets of features based on their partial correspondences.

[Previous work: Indyk & Thaper, Bartal, Charikar, Agarwal & Varadarajan, ...]

Kristen Grauman

### Pyramid match: main idea



Feature space partitions serve to "match" the local descriptors within successively wider regions.

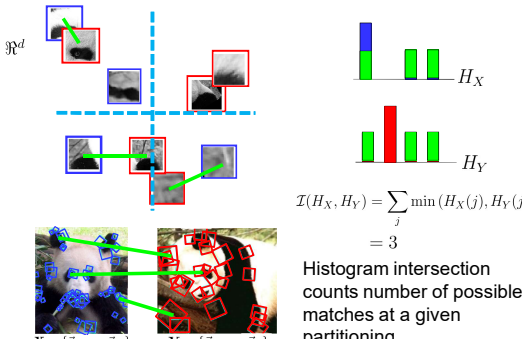
descriptor

face

$X = \{\vec{x}_1, \dots, \vec{x}_m\}$      $Y = \{\vec{y}_1, \dots, \vec{y}_n\}$

Kristen Grauman

### Pyramid match: main idea



$\mathcal{I}(H_X, H_Y) = \sum_j \min(H_X(j), H_Y(j))$   
 $= 3$

Histogram intersection counts number of possible matches at a given partitioning.

$X = \{\vec{x}_1, \dots, \vec{x}_m\}$      $Y = \{\vec{y}_1, \dots, \vec{y}_n\}$

Kristen Grauman

### Pyramid match kernel

$$K_{\Delta}(X, Y) = \sum_{i=0}^L 2^{-i} \left[ \mathcal{I}(H_X^{(i)}, H_Y^{(i)}) - \mathcal{I}(H_X^{(i-1)}, H_Y^{(i-1)}) \right]$$

measures difficulty of a match at level  $i$ 
number of newly matched pairs at level  $i$

- For similarity, weights inversely proportional to bin size (or may be learned)
- Normalize these kernel values to avoid favoring large sets

[Grauman & Darrell, ICCV 2005]

### Pyramid match kernel

**Optimal match:  $O(m^3)$**   
**Pyramid match:  $O(mL)$**

$X = \{\vec{x}_1, \dots, \vec{x}_m\}$      $Y = \{\vec{y}_1, \dots, \vec{y}_n\}$

Kristen Grauman

### Unordered sets of local features: No spatial layout preserved!

Too much?
Too little?

### Spatial pyramid match

- Make a pyramid of bag-of-words histograms.
- Provides some loose (global) spatial layout information

[Lazebnik, Schmid & Ponce, CVPR 2006]

### Spatial pyramid match

- Make a pyramid of bag-of-words histograms.
- Provides some loose (global) spatial layout information

$$K^L(X, Y) = \sum_{m=1}^M \kappa^L(X_m, Y_m)$$

Sum over PMKs computed in *image coordinate space*, one per word.

[Lazebnik, Schmid & Ponce, CVPR 2006]

### Spatial pyramid match

- Can capture **scene** categories well---texture-like patterns but with some variability in the positions of all the local

### Spatial pyramid match

- Can capture **scene** categories well---texture-like patterns but with some variability in the positions of all the local pieces.
- Sensitive to global shifts of the view

Confusion table

### SVMs: Pros and cons

- **Pros**
  - Kernel-based framework is very powerful, flexible
  - Often a sparse set of support vectors – compact at test time
  - Work very well in practice, even with very small training sample sizes
- **Cons**
  - No "direct" multi-class SVM, must combine two-class SVMs
  - Can be tricky to select best kernel function for a problem
  - Computation, memory
    - During training time, must compute matrix of kernel values for every pair of examples
    - Learning can take a very long time for large-scale problems

Adapted from Leo Leisner

### Recall: Evolution of methods

- Hand-crafted models
- 3D geometry
- Hypothesize and align
- Hand-crafted features
- Learned models
- Data-driven
- **"End-to-end" learning of features and models\*,\*\***

→

### Traditional Image Categorization: Training phase

Slide credit: lia-Bin Huang

### Traditional Image Categorization: Testing phase

Prediction  
**Outdoor**

Slide credit: lia-Bin Huang

### Learning a Hierarchy of Feature Extractors

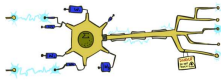
- Each layer of hierarchy extracts features from output of previous layer
- All the way from pixels → classifier
- Layers have the (nearly) same structure

- Train all layers jointly

Slide: Rob Fergus

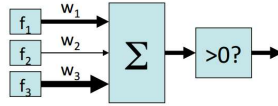
### Neuron: Linear Perceptron

- Inputs are **feature values**
- Each feature has a **weight**
- Sum is the **activation**



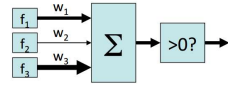
$$\text{activation}_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$

- If the activation is:
  - Positive, output +1
  - Negative, output -1



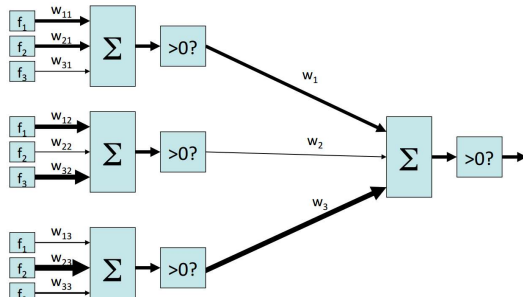
Slide credit: Pieter Abbeel and Dan Klein

### Two-layer perceptron network



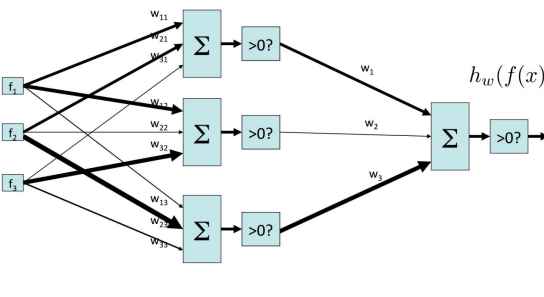
Slide credit: Pieter Abbeel and Dan Klein

### Two-layer perceptron network



Slide credit: Pieter Abbeel and Dan Klein

### Two-layer perceptron network



Slide credit: Pieter Abbeel and Dan Klein

### Convolutional Neural Networks (CNN, ConvNet, DCN)

- CNN = a multi-layer neural network with
  - Local connectivity:**
    - Neurons in a layer are only connected to a small region of the layer before it
  - Share weight parameters across spatial positions:**
    - Learning shift-invariant filter kernels

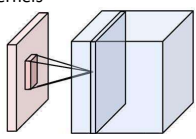
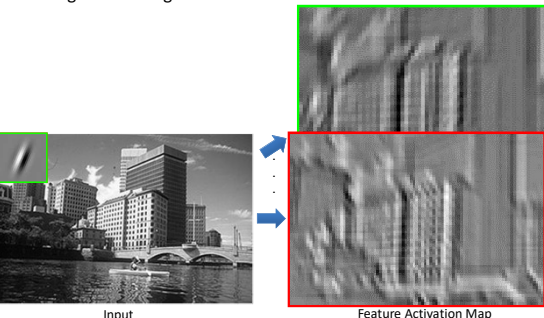


Image credit: A. Karpathy

Jia-Bin Huang and Derek Hoiem, UIUC

### What is a Convolution?

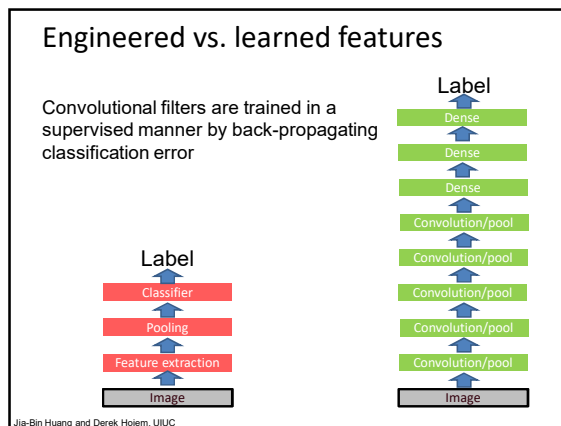
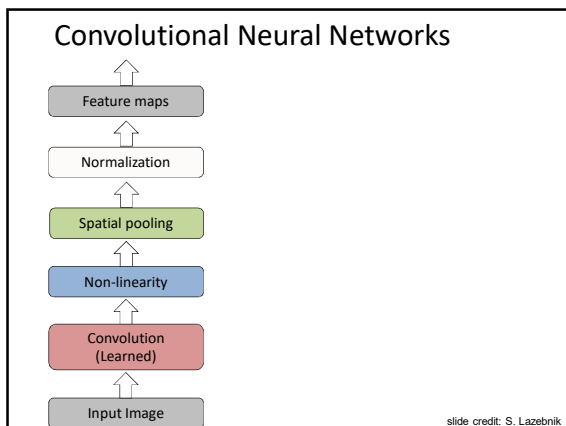
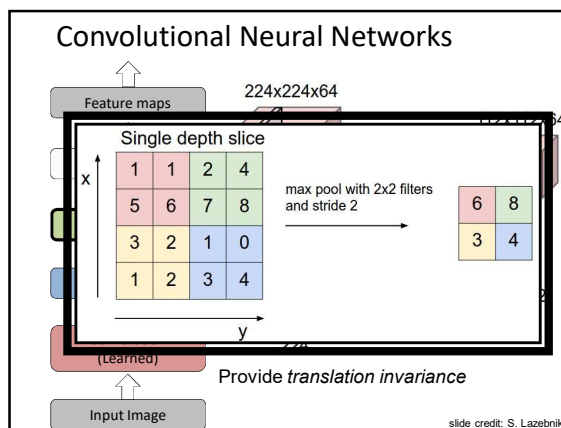
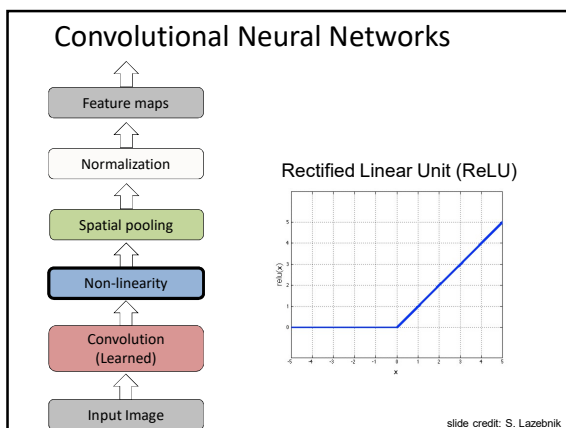
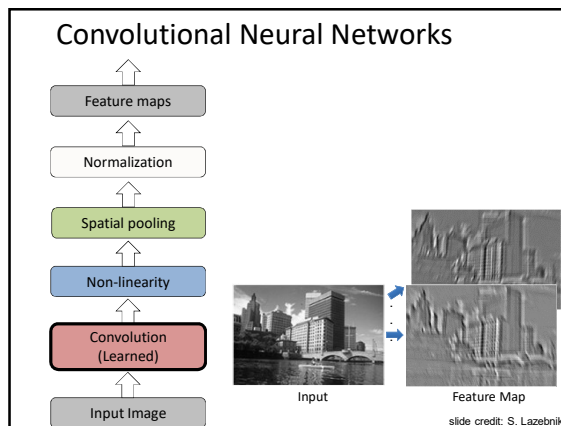
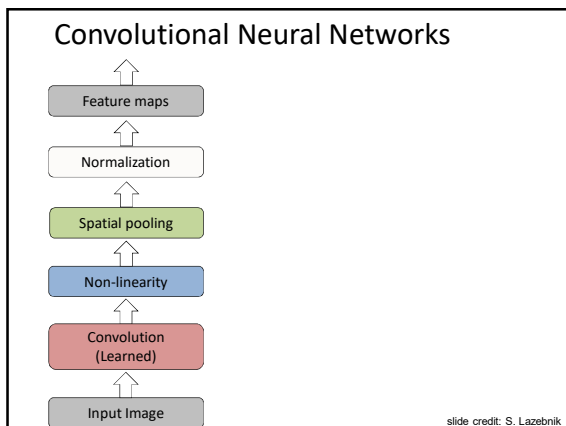
- Weighted moving sum



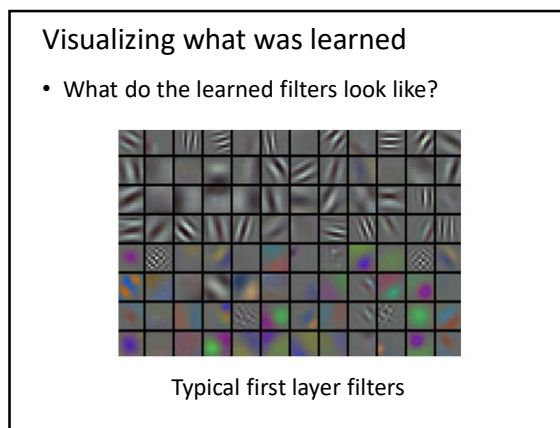
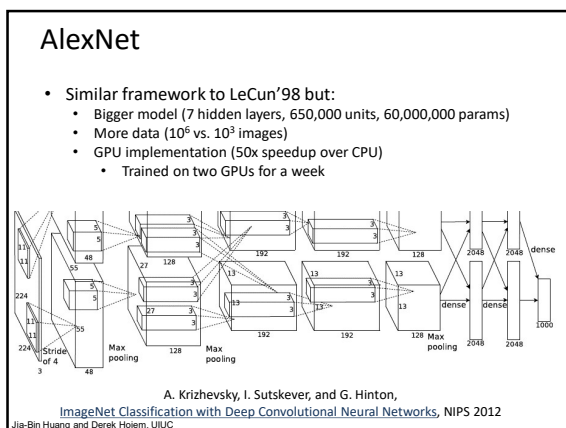
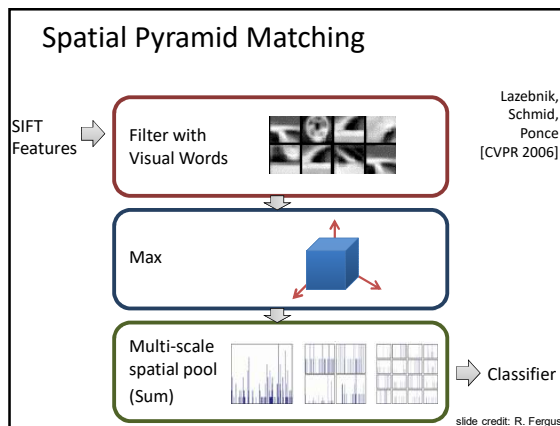
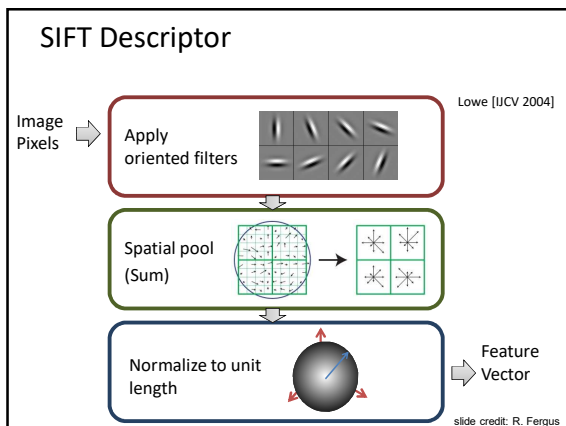
Input

Feature Activation Map

slide credit: S. Lazebnik







Wired

Google's Artificial Brain Learns to Find Cat Videos

SHARE

778

TWEET

COMMENT

EMAIL

BY LIAT CLARK, *Wired UK*

When computer scientists at Google's mysterious X lab built a neural network of 16,000 computer processors with one billion connections and let it browse YouTube, it did what

<https://www.wired.com/2012/06/google-x-neural-network/>

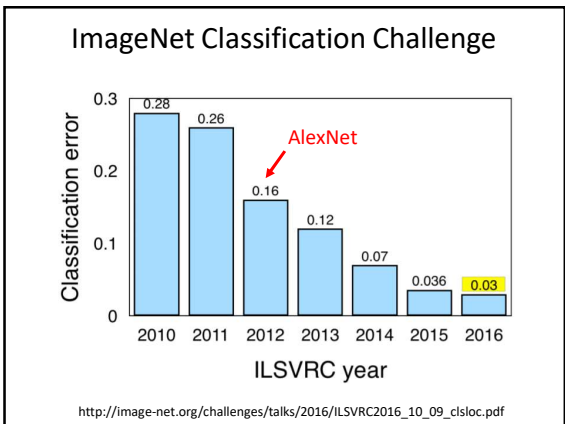
### Application: ImageNet

- ~14 million labeled images, 20k classes
- Images gathered from Internet
- Human labels via Amazon Turk

[Deng et al. CVPR 2009]

<https://sites.google.com/site/deeplearningcvpr2014>

Slide: R. Fergus



### Industry Deployment

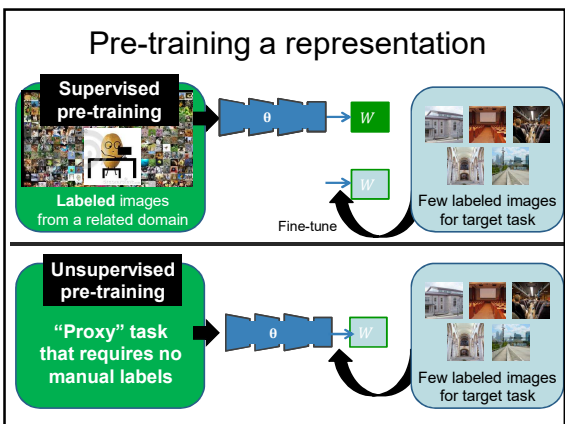
- Used in Facebook, Google, Microsoft
- Image Recognition, Speech Recognition, ...
- Fast at test time

Taigman et al. DeepFace: Closing the Gap to Human-Level Performance in Face Verification, CVPR'14

Slide: R. Fergus

- ### Beyond classification
- Detection
  - Segmentation
  - Regression
  - Pose estimation
  - Matching patches
  - Synthesis
- and many more...
- Jia-Bin Huang and Derek Hoiem, UIUC

- ### Recap
- Neural networks / multi-layer perceptrons
    - View of neural networks as learning hierarchy of features
  - Convolutional neural networks
    - Architecture of network accounts for image structure
    - “End-to-end” recognition from pixels
    - Together with big (labeled) data and lots of computation → major success on benchmarks, image classification and beyond



- ### New forms of self-supervision
- What can be our “proxy” or “pretext” task?
  - *Temporal coherence in video*
    - Mobahi et al. 2009, Wang & Gupta 2015, Wang et al. 2016, Gao et al. 2016, ...
  - *Audio channel – ambient sounds*
    - Owens et al. 2016, Arandjelovic & Zisserman 2017
  - *Ego-motion*
    - Jayaraman et al. 2015, Agrawal et al. 2015
  - *Spatial context, patch layout*
    - Doersch et al. 2015, Noroozi & Favaro 2016
  - *In-painting missing pixels*
    - Pathak et al. 2016
  - *Colorization*
    - Larsson et al. 2016, Zheng et al. 2016
  - *Temporal order of frames*
    - Misra et al. 2016

### Evaluation of self-supervised rep

How to test quality of unsupervised pre-training?


Comparisons against

- Equally supervised, but without unsup pretrain
- Fully supervised pre-training (ImageNet)
- Same network with random weights
- Counting "object-selective units" (Owens et al.)

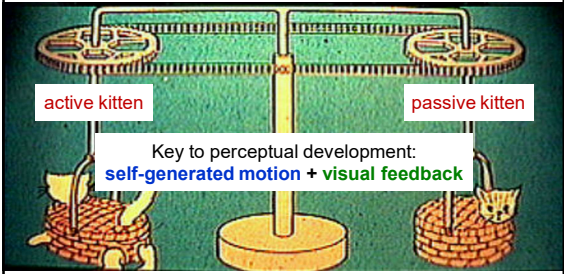
Raw representation, +/- fine-tuning to a task

### (Ego)motion for self-supervision

Dinesh Jayaraman and Kristen Grauman  
 Department of Computer Science  
 University of Texas at Austin



### The kitten carousel experiment [Held & Hein, 1963]




active kitten

passive kitten

Key to perceptual development:  
 self-generated motion + visual feedback


### Big picture goal: Embodied vision

**Status quo:**  
 Learn from "disembodied" bag of labeled snapshots.

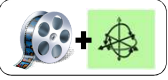



↓

**Goal:**  
 Learn in the context of **acting** and **moving** in the world.

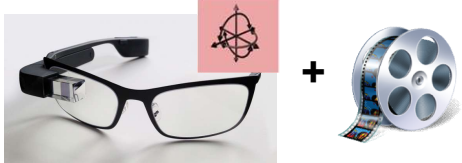


### Two formulations

1. Learning representations tied to ego-motion 
2. Learning representations from unlabeled video 

### Our idea: Ego-motion ↔ vision

**Goal:** Teach computer vision system the connection: "how I move" ↔ "how my visual surroundings change"




Ego-motion motor signals + Unlabeled video

[Jayaraman & Grauman, ICCV 2015]

**Our idea: Ego-motion ↔ vision**

**Goal:** Teach computer vision system the connection: "how I move" ↔ "how my visual surroundings change"

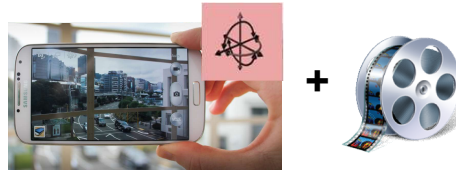


Ego-motion motor signals + Unlabeled video

[Jayaraman & Grauman, ICCV 2015]

**Our idea: Ego-motion ↔ vision**

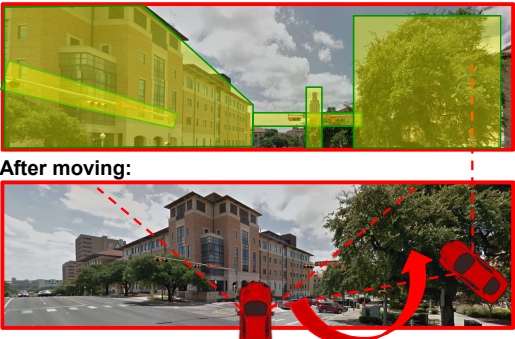
**Goal:** Teach computer vision system the connection: "how I move" ↔ "how my visual surroundings change"



Ego-motion motor signals + Unlabeled video

[Jayaraman & Grauman, ICCV 2015]

**Ego-motion ↔ vision: view prediction**



After moving:

**Approach idea: Ego-motion equivariance**

**Invariant features:** unresponsive to some classes of transformations

$$z(gx) \approx z(x)$$

- Simard et al, Tech Report, '98
- Wiskott et al, Neural Comp '02
- Hadsell et al, CVPR '06
- Mobahi et al, ICML '09
- Zou et al, NIPS '12
- Sohn et al, ICML '12
- Cadieu et al, Neural Comp '12
- Goroshin et al, ICCV '15
- Lies et al, PLoS computation biology '14
- ...

**Approach idea: Ego-motion equivariance**

**Invariant features:** unresponsive to some classes of transformations

$$z(gx) \approx z(x)$$

**Equivariant features:** *predictably* responsive to some classes of transformations, through simple mappings (e.g., linear)

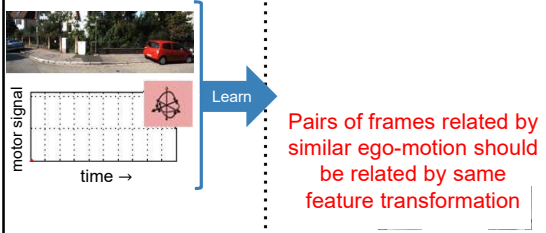
$$z(gx) \approx M_g z(x)$$

"equivariance map"

Invariance *discards* information; equivariance *organizes* it.

**Approach idea: Ego-motion equivariance**

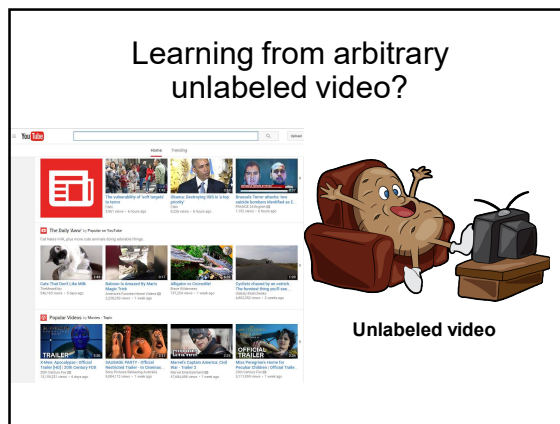
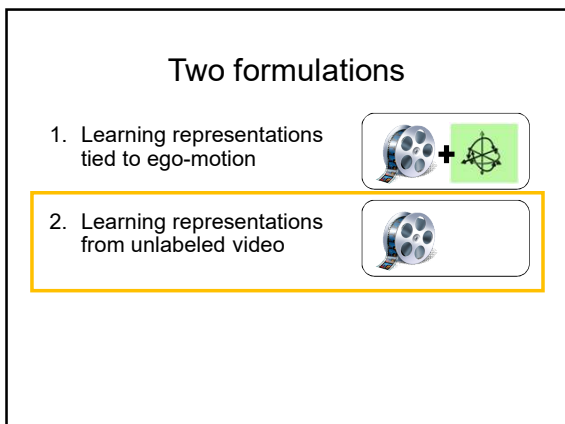
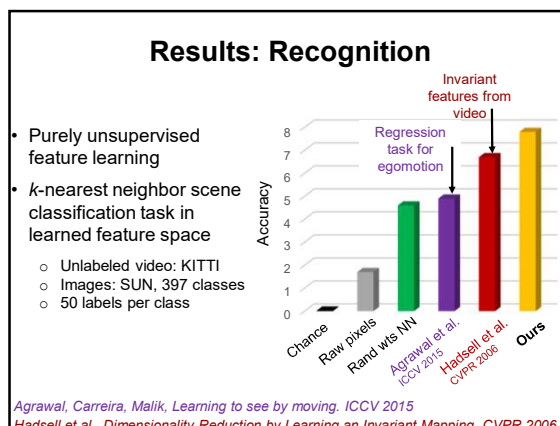
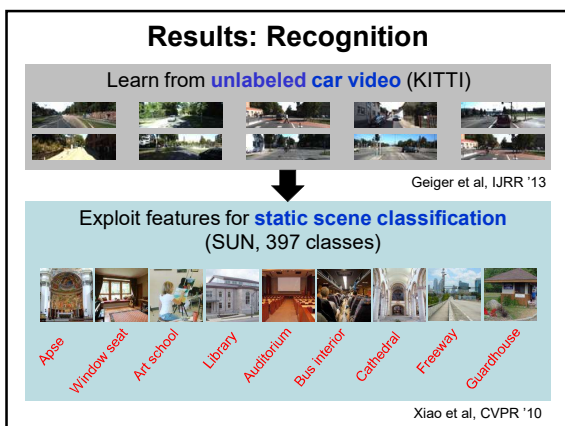
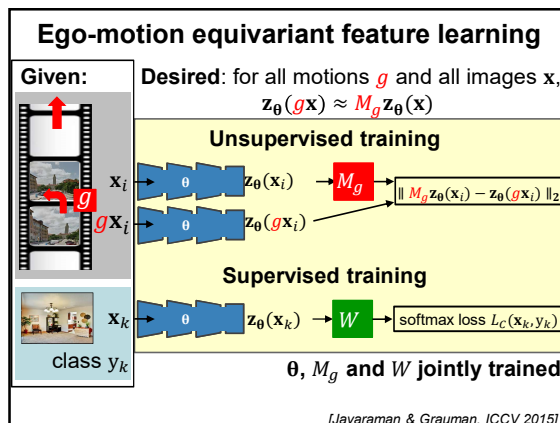
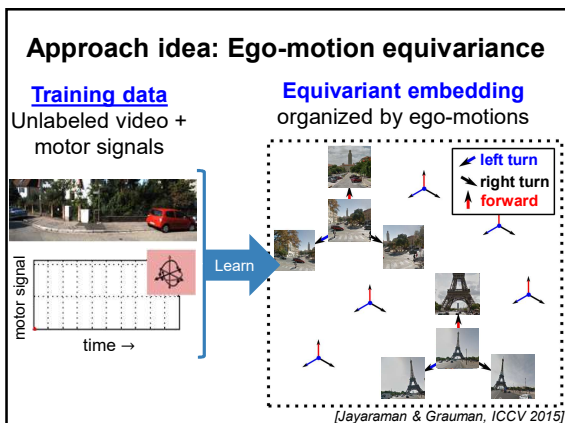
**Training data**  
Unlabeled video + motor signals



**Equivariant embedding**  
organized by ego-motions

Pairs of frames related by similar ego-motion should be related by same feature transformation

[Jayaraman & Grauman, ICCV 2015]



### Background: Slow feature analysis

[Wiskott & Sejnowski, 2002]

Find functions  $g(x)$  that map

quickly varying input signal  $x(t)$  → slowly varying features  $y(t)$

Figure: Laurenz Wiskott, <http://www.scholarpedia.org/article/File:SlowFeatureAnalysis-OptimizationProblem.png>

### Background: Slow feature analysis

[Wiskott & Sejnowski, 2002]

Find functions  $g(x)$  that map

quickly varying input signal  $x(t)$  → slowly varying features  $y(t)$

Figure: Laurenz Wiskott, <http://www.scholarpedia.org/article/File:SlowFeatureAnalysis-OptimizationProblem.png>

### Background: Slow feature analysis

[Wiskott & Sejnowski, 2002]

- Existing work exploits "slowness" as temporal coherence in video → learn invariant representation

[Hadsell et al. 2006; Mobahi et al. 2009; Bergstra & Bengio 2009; Goroshin et al. 2013; Wang & Gupta 2015,...]

$z(a) \approx z(b)$   
in learned embedding

- Fails to capture how visual content changes over time

### Our idea: Steady feature analysis

- Higher order temporal coherence in video → learn equivariant representation

Second order slowness operates on frame triplets:

$$z(b) - z(a) \approx z(c) - z(b)$$

in learned embedding

[Jayaraman & Grauman, CVPR 2016]

### Our idea: Steady feature analysis

Equivariance  $\approx$  "steadily" varying frame features!  
 $d^2 z_\theta(xt)/dt^2 \approx 0$

[Jayaraman & Grauman, CVPR 2016]

### Datasets

<b>Unlabeled video</b>	→	<b>Target task (few labels)</b>
Human Motion Database (HMDB)	→	PASCAL 10 Actions
KITTI Video	→	SUN 397 Scenes
NORB	→	NORB 25 Objects

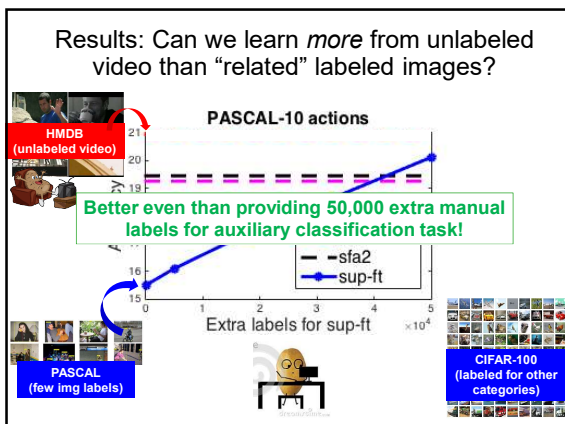
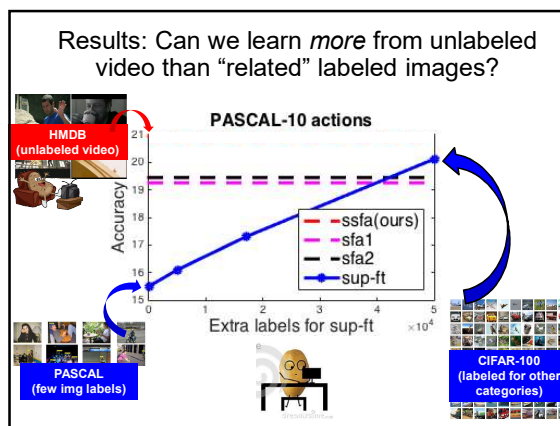
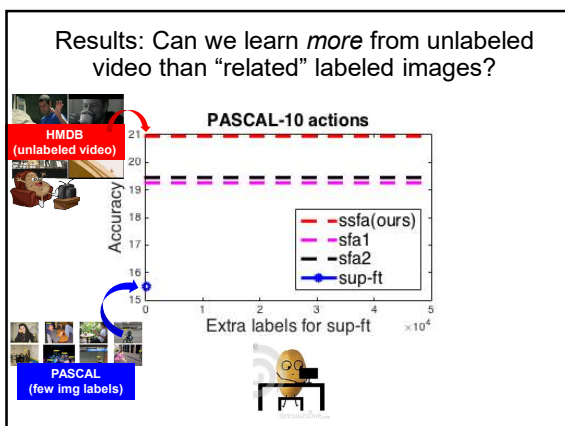
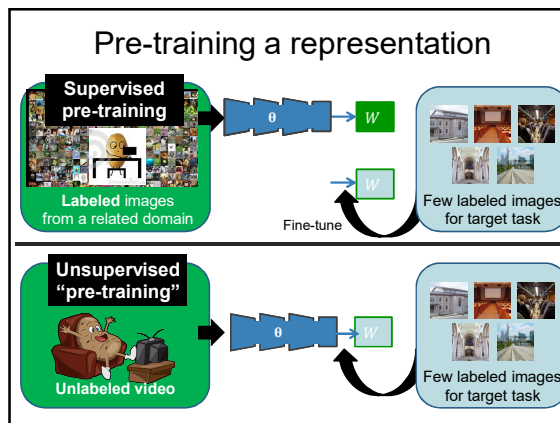
32 x 32 images or 96 x 96 images

### Results: Steady feature analysis

Task type →	Objects	Scenes	Actions	
Datasets →	NORB → NORB	KITTI → SUN	HMDB → PASCAL-10	
Methods ↓	[25 cls]	[397 cls]	[397 cls, top-10]	[10 cls]
random	4.00	0.25	2.52	10.00
UNREG	24.64 ± 0.85	0.70 ± 0.12	6.10 ± 0.67	15.34 ± 0.28
SFA-1 [30]*	37.57 ± 0.85	1.21 ± 0.14	8.24 ± 0.25	19.26 ± 0.45
SFA-2 [14]**	39.23 ± 0.94	1.02 ± 0.12	6.78 ± 0.32	19.04 ± 0.24
<b>SSFA (ours)</b>	<b>42.83 ± 0.33</b>	<b>1.65 ± 0.04</b>	<b>9.19 ± 0.10</b>	<b>20.95 ± 0.13</b>

Multi-class recognition accuracy

\*Hadsell et al., Dimensionality Reduction by Learning an Invariant Mapping, CVPR'06  
 \*\*Mobahi et al., Deep Learning from Temporal Coherence in Video, ICML'09



- ### Summary
- Visual learning benefits from
    - context of action and motion in the world
    - continuous self-acquired feedback
  - New ideas:
    - "Embodied" feature learning using both visual and motor signals
    - Feature learning from unlabeled video via higher order temporal coherence

## Papers

- **Learning Image Representations Tied to Ego-Motion.** D. Jayaraman and K. Grauman. In Proceedings of the IEEE International Conference on Computer Vision (ICCV), Santiago, Chile, Dec 2015.
- **Slow and Steady Feature Analysis: Higher Order Temporal Coherence in Video.** D. Jayaraman and K. Grauman. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Las Vegas, June 2016.