# designing and learning visual representations

stefano soatto

ucla

november 13, 2015

- **visual representation: definition**
  - tradeoffs

- **special (trivial) case: local descriptors**
  - the unreasonable effectiveness of sift

- **beyond local: modeling intrinsic variability**
  - conjectures

- **embedding the representation in the scene**
  - scene topology, gravity

# representation

$$\mathbb{I}(gT; \phi(x^t)) = \mathbb{H}(gT) - \mathbb{H}(gT|\phi(x^t)) \quad \forall \, g \in G$$

$$\underbrace{\max(0, \mathcal{G}(u,v;\sigma,\alpha) * x(u,v))}_{\text{ReLu}}$$

$\phi$ $x^t$ $\theta$ $T$ $g$

- a function of the data that is useful for a task ...

regardless of nuisance factors affecting (future) data

$$\mathcal{H}(y) = \mathbb{H}(y) - \mathbb{H}(\phi_{x^t,G}(y))$$

- data $\quad x^t = \{x_1, \ldots, x_t\}$

  - images

    $y$

- task $\quad \theta$

  $\theta$ $\qquad$ $g$
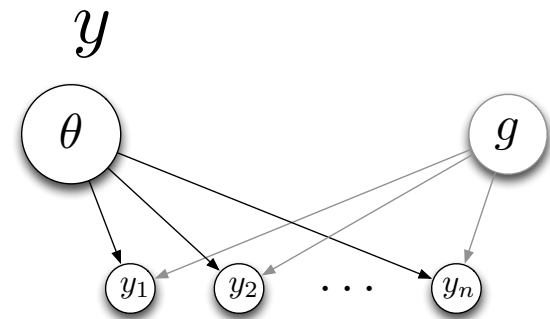
  $y_1$ $y_2$ $\cdots$ $y_n$

  - decision or control actions on the scene portrayed by the images

- nuisance factors $\quad g$

  - viewpoint, illumination, partial occlusion, sensor characteristics

- useful

  $$\theta_k \sim p(\theta|k)$$

  - "informative"

  $$x_k \sim p_{\theta_k}(x)$$

$$= \int p_\theta(y|x_k, g_k) dP(g_k|k) \doteq \int p_{\theta_k, g_k}(y) dP(g_k|k) \qquad (9)$$

# optimal representation

- 'most informative' function of the data, for a task: sufficient statistic

- 'most compressed': minimal sufficient statistic

- 'most insensitive' to nuisance factors affecting future data: minimal sufficient invariant statistic

# (fake) bad news

- 'invariance' cannot be attained, so settle for 'approximate invariance'

- invariance trades off information (discriminative power)

# a few facts

- The likelihood function $L(\theta) \doteq p_\theta(x)$ is a minimal sufficient statistic, even if $\theta$ is infinite-dimensional (Bahadur, 1954).
- Nuisance variability can be <span style="color:red">marginalized</span>: $p_\theta(x|G) \doteq \int_G p_\theta(x|g)dP(g)$

  but invariant only if marginalization is wrt the base measure and in general *not maximal.*
- <span style="color:red">Profiling</span> (max-out): $p_{\theta,G}(x) \doteq \sup_{g \in G} p_{\theta,g}(x)$

  yields *a maximal invariant*, but (non-convex) search at test time.
- <span style="color:red">(Down)-sampling</span> the profile likelihood introduces <span style="color:red">aliasing</span> phenomena (extrema that do not exist before downsampling/reconstruction)
- <span style="color:red">Anti-aliasing = pooling = local marginalization</span>

$$\hat{p}_{\theta,G}(y) = \max_i \hat{p}_{\theta,g_i}(y) = \max_i \int_G p_{\theta,g_i}(gy)w(g)d\mu(g)$$

$p_\theta$

the sample-orbit antialiased likelihood is a minimal sufficient invariant statistic: optimal representation
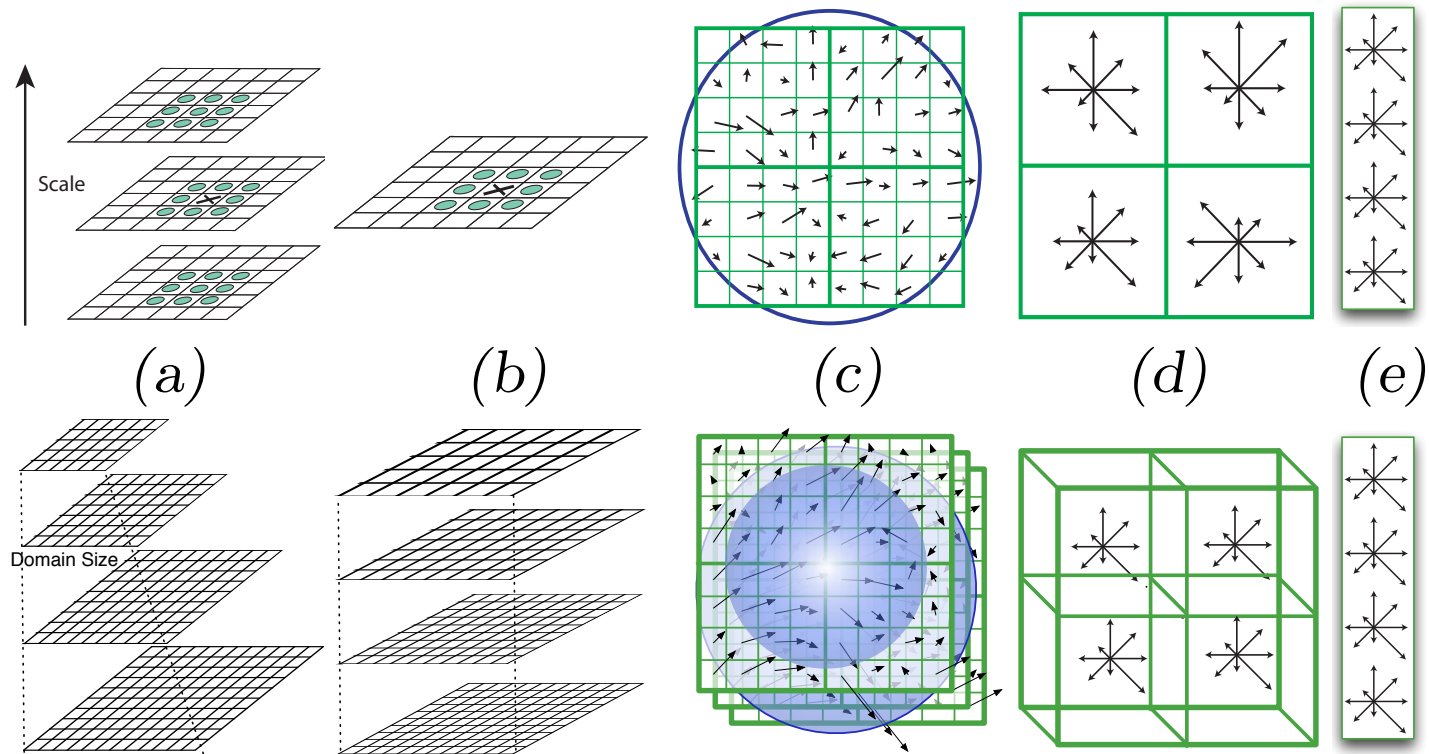
# simple example

- training set: a single image

- task: binary classification (correspondence)

- nuisances: planar similarities

- optimal representation (closed form):

$$p_x(y|\mathcal{H}) \doteq p(\angle\nabla y | \nabla x) = \frac{1}{\sqrt{2\pi\epsilon^2}} \exp\left(-\frac{1}{2\epsilon^2}\sin^2(\angle\nabla y - \angle\nabla x)\|\nabla x\|^2\right) M$$
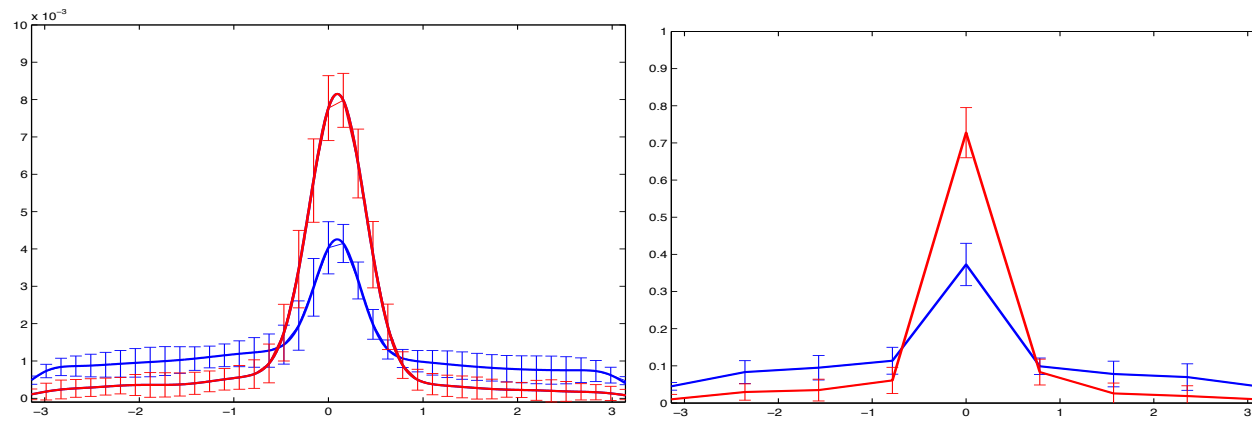
$$M = \frac{\epsilon e^{-\frac{(m)^2}{2\epsilon^2}}}{\sqrt{2\pi}} + m - m\Psi\left(-\frac{m}{\epsilon}\right)$$

# SIFT revisited

$$h_{\mathrm{SIFT}}(\theta|I) = \int \kappa_\epsilon\left(\theta - \angle\nabla I(y)\right)\kappa_\sigma(y - x)\|\nabla I(y)\|dy$$



(a)          (b)          (c)          (d)          (e)

$$h_{DSP}(\theta|I) = \iint \kappa_\epsilon\left(\theta - \angle\nabla I(y)\right)\kappa_\sigma(y - x)\|\nabla I(y)\|dP(\sigma)dy$$

w/ j. dong

CLAMPING

# multi-view descriptors

- MV-HOG $\quad \phi_{\mathbf{z}}^t(\theta|\mathbf{y}^t) \doteq \dfrac{1}{\tau} \displaystyle\sum_{\tau=1}^{t} \int \mathcal{N}_{\mathbb{S}^1}(\theta - \angle g\mathbf{y}(\tau)) \|\nabla\mathbf{y}\| dP(g)$

- R-HOG $\quad \phi_{\hat{\mathbf{z}}}^t(\theta) \doteq \displaystyle\int_{SO(3)\times\mathbb{R}} \mathcal{N}_{\mathbb{S}^1}(\theta - \angle g\hat{\mathbf{z}}) dP_{SO(3)}(g) \|\nabla\hat{\mathbf{z}}\| d\mu$

**w/ j. dong, j. hernandez, d. davis, j. balzer**

# where are we?

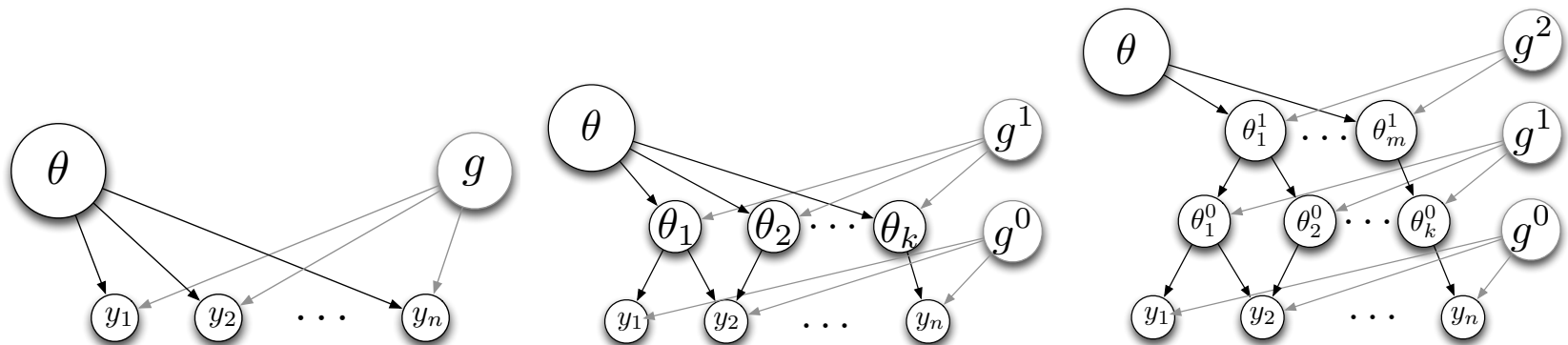- structure in the representation?

- topology? geometry?

$$p(\theta|k)$$

$$\|(v)\|dv$$

# more realistically

$$\simeq \max_{i,V} \int_{G^M} \prod_{j\in V} \hat{p}_{\theta_k,g_k g_i}(y_{|g_j \mathcal{B}_0}) dP(g_1^{-1}g_k,\ldots,g_M^{-1}g_k|k)$$

$$= \max_{i,V} \int_{G^M} \prod_{j\in V} \hat{p}_{\theta_k,g_k g_i}(g_j y_{|\mathcal{B}_0}) dP(g_1^{-1}g_k,\ldots,g_M^{-1}g_k|k)$$

$(\sigma)$  ● intra-class variability (sepa ... principle) $\qquad x_k$

$$= \max_V \int_{G^M} \prod_{j\in V} \hat{p}_{\theta_k,g_{k_j}}(y) dP_G(g_{k_1},\ldots,g_{k_M}|k)$$

$$p_\theta(y|k) = \int p_\theta(y|x_k,g_k) dP(g_k|k) \doteq \int p_{\theta_k,g_k}(y) dP(g_k|k)$$

$x = \theta$  $\{V, \{g_{k_j}\}_{j\in V}, \{y_{|g_{k_j}\mathcal{B}_0}\}_{j\in V}\}$. $V$ determines the visible components, $\{g_{k_j}\}$

● occlusion (combinatorics) their photometric appearance.

their geometric configuration, and $\{y_{|g_{k_j}\mathcal{B}_0}\}$

$g_{k_j}$ the restriction of the group action $g_k$ on the domain of the receptive field $g_j$

$g$  $\qquad p_{x,G}(y) = p(\angle\nabla y|x)\|\nabla x\|$   $\qquad \hat{p}_{\theta,G,\hat{V}}(y|k) = \max_{i,V\in\mathcal{P}(D)} \int_{\text{diff}(D)} \prod_{j\in V} \hat{p}_{\theta_k,g_i}(y_j|x_k,g_k) dP(g_k|k)$

$= \prod_i \mathcal{N}_{g_i}(\alpha_i - \angle\nabla x_i;\epsilon_i)\|\nabla x_i\|$

$= \int_{SE(2)\times\mathbb{R}^+} \mathcal{N}_\epsilon(\alpha - \angle\nabla x(v))\mathcal{N}_\sigma(u-v)\|\nabla x(v)\|dv$

$\simeq \max_{i,V}\int_{G^M}\prod_{j\in V}\hat{p}_{\theta_k,g_k g_i}(y_{|g_j \mathcal{B}_0})dP(g_1^{-1}g_k,\ldots,g_M^{-1}g_k|k)$

$$\underbrace{\max(0,\mathcal{G}(u,v;\sigma,\alpha)*x(u,v))}_{\text{ReLu}} = \max_{i,V}\int_{G^M}\prod_{j\in V}\underbrace{\mathcal{N}(u-\tilde{u},v-\tilde{v};\sigma)\kappa_+(\angle\nabla x(\tilde{u},\tilde{v}),\alpha)\|\nabla x(\tilde{u},\tilde{v})\|d\tilde{u}d\tilde{v}}_{\text{SIFT}}$$

$- \angle\nabla x_{\tau j};\epsilon_\alpha)\mathcal{N}_\sigma(i-j)\mathcal{E}_s(y)d\mu(j)dP(\sigma)$  ReLu

$\hat{p}_{\theta_k,g_k g_i}(g_j y_{|\mathcal{B}_0})dP(g_1^{-1}g_k,\ldots,g_M^{-1}g_k|k)$

$$= \max_V \int_{G^M} \prod_{j\in V} \hat{p}_{\theta_k,g_{k_j}}(y) dP_G(g_{k_1},\ldots,g_{k_M}|k)$$

● architecture? $\{V,\{g_{k_j}\}_{j\in V},\{y_{|g_{k_j}\mathcal{B}_0}\}_{j\in V}\}$. $V$ determines the visible components, $\{g_{k_j}\}$
their geometric configuration, and $\{y_{|g_{k_j}\mathcal{B}_0}\}$ their photometric appearance.

$x^t$ $\qquad \theta$

$gT) - \mathbb{H}(gT|\phi(x^t))$

$) - \mathbb{H}(\phi_{x^t,G}(y))$

# still missing

- scene topology (detachable objects)

- global referencing (gravity)

- extension to tasks other than detection

# summary

- definition, analytical characterization of an ideal visual representation:

- simple case where inference is tractable

- conjectures on extensions to include intra-class variability: relation to cnn's

- representation of the scene, not the image

- extension to more general (control, decision) tasks

- support "query system" on the scene