# Microsoft Research

Each year Microsoft Research hosts hundreds of influential speakers from around the world including leading scientists, renowned experts in technology, book authors, and leading academics, and makes videos of these lectures freely available.

# On Decomposing the Proximal Map

Yao-Liang Yu

University of Alberta

December 6, 2013

# Regularized loss minimization

Generic form for many ML problems:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \ell(\mathbf{w}) + f(\mathbf{w})$$

- $\ell$ is the loss function;
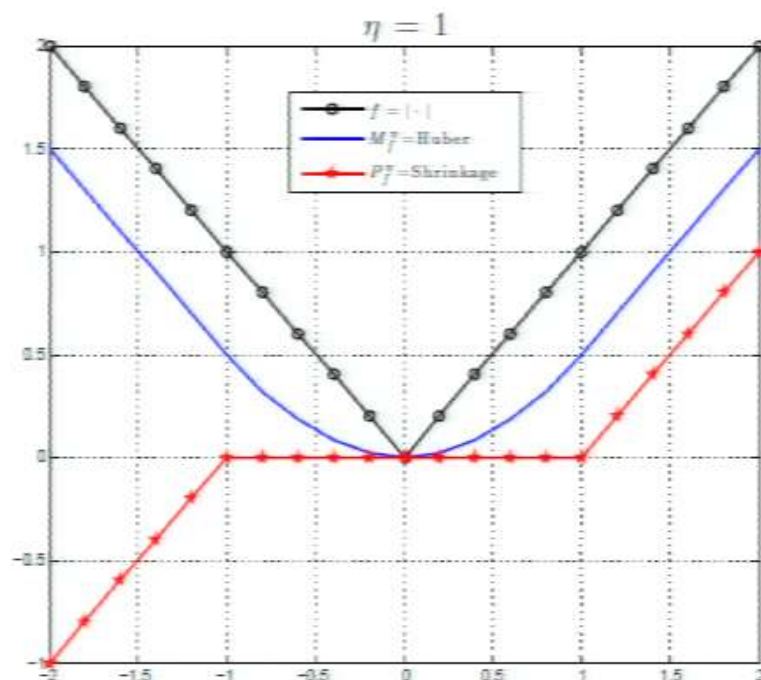- $f$ is the regularizer, usually a (semi)norm;

Special interest:

- sparsity;
- computational efficiency.

# Moreau envelop and proximal map

**Definition (Moreau'65)**

$$M_f(\mathbf{y}) = \min_{\mathbf{w}} \tfrac{1}{2}\|\mathbf{w} - \mathbf{y}\|^2 + f(\mathbf{w})$$

$$P_f(\mathbf{y}) = \operatorname*{argmin}_{\mathbf{w}} \tfrac{1}{2}\|\mathbf{w} - \mathbf{y}\|^2 + f(\mathbf{w})$$

$\eta = 1$

- $f = |\cdot|$
- $M_f^\eta = \text{Huber}$
- $P_f^\eta = \text{Shrinkage}$

# Proximal gradient (Fukushima & Mine'81)

$$\min_{\mathbf{w} \in \mathbb{R}^d} \ell(\mathbf{w}) + f(\mathbf{w})$$

> ❶ $\mathbf{y}_t = \mathbf{w}_t - \eta \nabla \ell(\mathbf{w}_t);$
>
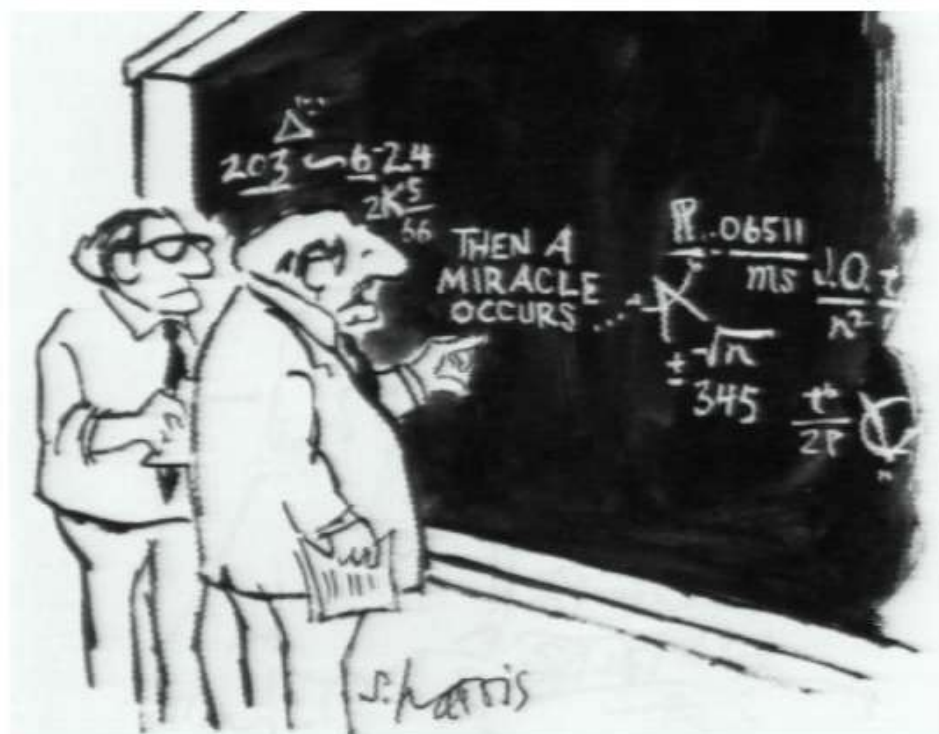> ❷ $\mathbf{w}_{t+1} = \mathsf{P}_{\eta f}(\mathbf{y}_t).$

For $f = \|\cdot\|_1$, obtain the shrinkage operator

$$[\mathsf{P}_{\|\cdot\|_1}(\mathbf{y})]_i = \mathrm{sign}(y_i)(|y_i| - 1)_+.$$

- guaranteed convergence, can be accelerated;
- generalization of projected gradient: $f = \iota_C$;
- reveals the sparsity-inducing property.

Refs: Combettes & Wajs'05; Beck & Teboulle'09; Duchi & Singer'09; Nesterov'13; etc.

# Then A Miracle Occurs...



"I think you should be more explicit here in step two."

from *What's so Funny about Science?* by Sidney Harris (1977)

Step 2: $P_f(\mathbf{y}) = \underset{\mathbf{w}}{\mathrm{argmin}}\ \frac{1}{2}\|\mathbf{y} - \mathbf{w}\|^2 + f(\mathbf{w})$

# How to decompose?

- Typical structured sparse regularizers:

$$f(\mathbf{w}) = \sum_i f_i(\mathbf{w});$$

**Theorem (Parallel Sum)**

$$P_{f+g} = (P_{2f}^{-1} + P_{2g}^{-1})^{-1} \circ (2\mathsf{Id}).$$

- Not directly useful due to the inversion;
- Can numerically reduce to $P_f$ and $P_g$ (Combettes et al.'11);
- But a two-loop routine can be as slow as subgradient descent (Schmidt et.al'11; Villa et al.'13).

# How to decompose?

- Typical structured sparse regularizers:

$$f(\mathbf{w}) = \sum_i f_i(\mathbf{w});$$

## Theorem (Parallel Sum)

$$P_{f+g} = (P_{2f}^{-1} + P_{2g}^{-1})^{-1} \circ (2\mathrm{Id}).$$

- Not directly useful due to the inversion;
- Can numerically reduce to $P_f$ and $P_g$ (Combettes et al.'11);
- But a two-loop routine can be as slow as subgradient descent (Schmidt et.al'11; Villa et al.'13).

# How to decompose?

- Typical structured sparse regularizers:

$$f(\mathbf{w}) = \sum_i f_i(\mathbf{w});$$

## Theorem (Parallel Sum)

$$P_{f+g} = (P_{2f}^{-1} + P_{2g}^{-1})^{-1} \circ (2\mathsf{Id}).$$

- Not directly useful due to the inversion;
- Can numerically reduce to $P_f$ and $P_g$ (Combettes et al.'11);
- But a two-loop routine can be as slow as subgradient descent (Schmidt et.al'11; Villa et al.'13).

# How to decompose?

- Typical structured sparse regularizers:

$$f(\mathbf{w}) = \sum_i f_i(\mathbf{w});$$

## Theorem (Parallel Sum)

$$\mathsf{P}_{f+g} = (\mathsf{P}_{2f}^{-1} + \mathsf{P}_{2g}^{-1})^{-1} \circ (2\mathsf{Id}).$$

- Not directly useful due to the inversion;
- Can numerically reduce to $\mathsf{P}_f$ and $\mathsf{P}_g$ (Combettes et al.'11);
- But a two-loop routine can be as slow as subgradient descent (Schmidt et.al'11; Villa et al.'13).

# Two previous results

## Theorem (Friedman et al.'07)

$$P_{\|\cdot\|_1 + \|\cdot\|_{TV}} = P_{\|\cdot\|_1} \circ P_{\|\cdot\|_{TV}}, \quad \text{where} \quad \|\mathbf{w}\|_{TV} = \sum_{i=1}^{d-1} |w_i - w_{i+1}|.$$

## Theorem (Jenatton et al.'11)

*Assuming the groups $\{g_i\}$ form a laminar system ($g_i \cap g_j \in \{g_i, g_j, \emptyset\}$), then, if appropriately ordered,*

$$P_{\sum_{i=1}^{k} \|\cdot\|_{g_i}} = P_{\|\cdot\|_{g_1}} \circ \cdots \circ P_{\|\cdot\|_{g_k}},$$

*where $\|\cdot\|_{g_i}$ is the restriction of $l_p, p \in \{1, 2, \infty\}$ to the group $g_i$.*

## Generalization

$$P_{f+g} \overset{?}{=} P_f \circ P_g \overset{?}{=} P_g \circ P_f.$$

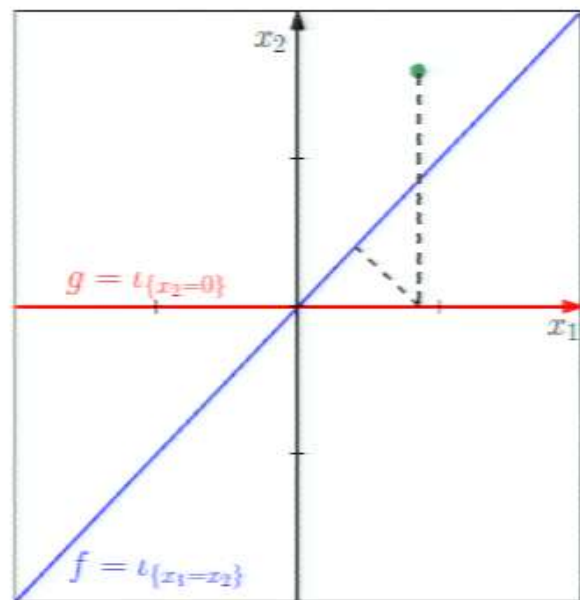But, is it even sensible?

# Bad news

## Theorem

On the real line, $\exists h$ such that $P_h = P_f \circ P_g$.

- Not necessarily $h = f + g$, though

## Example (A simple counterexample)

Consider $\mathbb{R}^2$, and let $f = \iota_{\{x_1 = x_2\}}, g = \iota_{\{x_2 = 0\}}$.



$$P_f = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}, \quad P_g = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}.$$

But $P_f \circ P_g = \begin{bmatrix} 0.5 & 0 \\ 0.5 & 0 \end{bmatrix}$

no $h$ such that $P_h = P_f \circ P_g$

# Nevertheless

- Can ask the decomposition to hold for many but not all cases.
- Manipulating the optimality conditions:

$$P_{f+g}(\mathbf{z}) = \mathrm{argmin}_{\mathbf{w}} \tfrac{1}{2} \|\mathbf{z} - \mathbf{w}\|^2 + (f + g)(\mathbf{w})$$
$$P_g(\mathbf{z}) = \mathrm{argmin}_{\mathbf{w}} \tfrac{1}{2} \|\mathbf{z} - \mathbf{w}\|^2 + g(\mathbf{w})$$
$$P_f(P_g(\mathbf{z})) = \mathrm{argmin}_{\mathbf{w}} \tfrac{1}{2} \|P_g(\mathbf{z}) - \mathbf{w}\|^2 + f(\mathbf{w}).$$

## Theorem

A sufficient condition for $P_{f-g}(\mathbf{z}) = P_f(P_g(\mathbf{z}))$ is

$$\forall \mathbf{y} \in \mathrm{dom}\, g, \ \partial g(P_f(\mathbf{y})) \supseteq \partial g(\mathbf{y}).$$

- Fails to be necessary at boundary points
- A special case appeared in a proof of (Zhou et al.'12)

# Nevertheless

- Can ask the decomposition to hold for many but not all cases.
- Manipulating the optimality conditions:

$$P_{f+g}(z) - z + \partial(f+g)(P_{f+g}(z)) \ni 0$$
$$P_g(z) - z + \partial g(P_g(z)) \ni 0$$
$$P_f(P_g(z)) - P_g(z) + \partial f(P_f(P_g(z))) \ni 0.$$

### Theorem

A sufficient condition for $P_{f+g}(z) = P_f(P_g(z))$ is

$$\forall\, y \in \operatorname{dom} g, \ \partial g(P_f(y)) \supseteq \partial g(y).$$

- Fails to be necessary at boundary points
- A special case appeared in a proof of (Zhou et al. '12)

# Nevertheless

- Can ask the decomposition to hold for many but not all cases.
- Manipulating the optimality conditions:

$$P_{f+g}(z) - z + \partial(f + g)(P_{f+g}(z)) \ni 0$$
$$P_f(P_g(z)) - z + \partial g(P_g(z)) + \partial f(P_f(P_g(z))) \ni 0.$$

## Theorem

A sufficient condition for $P_{f-g}(z) = P_f(P_g(z))$ is

$$\forall y \in \text{dom } g, \; \partial g(P_f(y)) \supseteq \partial g(y).$$

- Fails to be necessary at boundary points
- A special case appeared in a proof of (Zhou et al.'12)

# Nevertheless

- Can ask the decomposition to hold for many but not all cases.
- Manipulating the optimality conditions:

$$P_{f+g}(z) - z + \partial(f+g)(P_{f+g}(z)) \ni 0$$
$$P_f(P_g(z)) - z + \partial g(P_g(z)) + \partial f(P_f(P_g(z))) \ni 0.$$

---

**Theorem**

A sufficient condition for $P_{f+g}(z) = P_f(P_g(z))$ is

$$\forall\, y \in \operatorname{dom} g, \; \partial g(P_f(y)) \supseteq \partial g(y).$$

---

- Fails to be necessary at boundary points
- A special case appeared in a proof of (Zhou et al.'12)

# The rest is easy



- Find $f$ and $g$ that clinch our sufficient condition.

# Start with "trivialities"

---

**Theorem**

Fix $f$. $P_{f+g} = P_f \circ P_g$ for *all* $g$ if and only if

- $\dim(\mathcal{H}) \geq 2$; $f \equiv c$ or $f = \iota_{\{w\}} + c$ for some $c \in \mathbb{R}$ and $w \in \mathcal{H}$;
- $\dim(\mathcal{H}) = 1$ and $f = \iota_C + c$ for some $c \in \mathbb{R}$ and set $C$ that is closed and convex.

---

Asymmetry.

---

**Theorem**

Fix $g$. $P_{f+g} = P_f \circ P_g$ for *all* $f$ if and only if $g$ is continuous affine.

---

- Reassuring the impossibility to always have $P_{f+g} = P_f \circ P_g$;
- Still hope to get interesting results!

# Scaling Invariant $\Leftrightarrow$ Positive Homogeneous

$$\partial g(\mathsf{P}_f(\mathbf{y})) \supseteq \partial g(\mathbf{y})$$

$g$ positive homogeneous $\Leftrightarrow \forall \lambda > 0, \partial g(\lambda \mathbf{w}) = \partial g(\mathbf{w}) \Rightarrow \forall \mathbf{z}, \mathsf{P}_f(\mathbf{z}) \propto \mathbf{z}$

# Scaling Invariant $\Leftrightarrow$ Positive Homogeneous

$$\partial g(\mathsf{P}_f(y)) \supseteq \partial g(y)$$

$g$ positive homogeneous $\Leftrightarrow \forall \lambda > 0, \partial g(\lambda \mathbf{w}) = \partial g(\mathbf{w}) \Rightarrow \forall \mathbf{z}, \mathsf{P}_f(\mathbf{z}) \propto \mathbf{z}$

## Theorem

Fix $f$. The following are equivalent (provided $\dim(\mathcal{H}) \geq 2$):

i).

ii).

iii). For all $z \in \mathcal{H}$, $\mathsf{P}_f(z) = \lambda_z \cdot z$ for some $\lambda_z \in [0, 1]$;

iv).

# Scaling Invariant $\Leftrightarrow$ Positive Homogeneous

$$\partial g(P_f(\mathbf{y})) \supseteq \partial g(\mathbf{y})$$

$g$ positive homogeneous $\Leftrightarrow \forall \lambda > 0, \partial g(\lambda \mathbf{w}) = \partial g(\mathbf{w}) \Rightarrow \forall \mathbf{z}, P_f(\mathbf{z}) \propto \mathbf{z}$

## Theorem

Fix $f$. The following are equivalent (provided $\dim(\mathcal{H}) \geq 2$):

i). $f = h(\|\cdot\|)$ for some increasing function $h : \mathbb{R}_+ \to \mathbb{R} \cup \{\infty\}$;

ii). For all perpendicular $\mathbf{x} \perp \mathbf{y}$, $f(\mathbf{x} + \mathbf{y}) \geq f(\mathbf{y})$;

iii). For all $\mathbf{z} \in \mathcal{H}$, $P_f(\mathbf{z}) = \lambda_{\mathbf{z}} \cdot \mathbf{z}$ for some $\lambda_{\mathbf{z}} \in [0, 1]$;

iv). $\mathbf{0} \in \mathrm{dom} f$ and $P_{f+\kappa} = P_f \circ P_\kappa$ for all positive homogeneous $\kappa$.

If $\dim(\mathcal{H}) = 1$, only ii) $\implies$ i) ceases to hold.

# Scaling Invariant $\Leftrightarrow$ Positive Homogeneous

$$\partial g(\mathsf{P}_f(\mathbf{y})) \supseteq \partial g(\mathbf{y})$$

$g$ positive homogeneous $\Leftrightarrow \forall \lambda > 0, \partial g(\lambda \mathbf{w}) = \partial g(\mathbf{w}) \Rightarrow \forall \mathbf{z}, \mathsf{P}_f(\mathbf{z}) \propto \mathbf{z}$

## Theorem

Fix $f$. The following are equivalent (provided $\dim(\mathcal{H}) \geq 2$):

i). $f = h(\|\cdot\|)$ for some increasing function $h : \mathrm{R}_+ \to \mathrm{R} \cup \{\infty\}$;

ii). For all perpendicular $\mathbf{x} \perp \mathbf{y}$, $f(\mathbf{x} + \mathbf{y}) \geq f(\mathbf{y})$;

iii). For all $\mathbf{z} \in \mathcal{H}$, $\mathsf{P}_f(\mathbf{z}) = \lambda_{\mathbf{z}} \cdot \mathbf{z}$ for some $\lambda_{\mathbf{z}} \in [0, 1]$;

iv). $\mathbf{0} \in \mathrm{dom}\, f$ and $\mathsf{P}_{f+\kappa} = \mathsf{P}_f \circ \mathsf{P}_\kappa$ for *all* positive homogeneous $\kappa$.

If $\dim(\mathcal{H}) = 1$, only ii) $\implies$ i) ceases to hold.
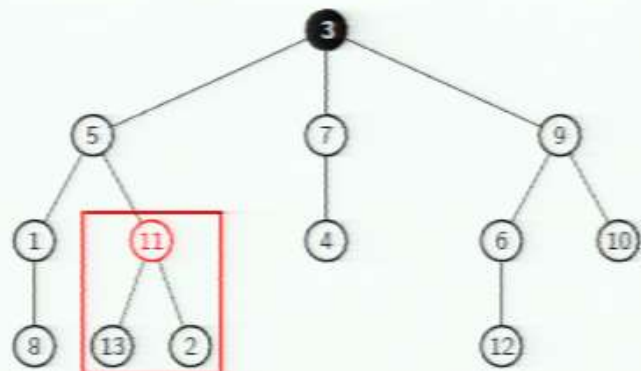
# Some Implications

## Theorem

*Fix $f$. The following are equivalent (provided $\dim(\mathcal{H}) \geq 2$):*

  i). $f = h(\|\cdot\|)$ *for some increasing function* $h : \mathbb{R}_+ \to \mathbb{R} \cup \{\infty\}$;

  ii). *For all perpendicular* $\mathbf{x} \perp \mathbf{y}$, $f(\mathbf{x} + \mathbf{y}) \geq f(\mathbf{y})$;

  iii). *For all* $z \in \mathcal{H}$, $P_f(z) = \lambda_z \cdot z$ *for some* $\lambda_z \in [0, 1]$;

  iv). $0 \in \operatorname{dom} f$ *and* $P_{f+\kappa} = P_f \circ P_\kappa$ *for all positive homogeneous $\kappa$.*

*If $\dim(\mathcal{H}) = 1$, only ii) $\implies$ i) ceases to hold.*

## i) $\Longleftrightarrow$ ii)

- Characterizing representer theorem (Dinuzzo & Schölkopf'12);
- Now we have more.

# Some Implications

## i) $\implies$ iv)

$$P_{\lambda\|\cdot\|^2+\kappa} = P_{\lambda\|\cdot\|^2} \circ P_\kappa = \tfrac{1}{\lambda+1}P_\kappa$$

- Double shrinkage;
- $\kappa = \|\cdot\|_1$: Elastic net (Zou & Hastie'05);
- Adding an $l_2$-ish regularizer, computationally, is free.

# Some Implications

## Theorem

Fix $f$. The following are equivalent (provided $\dim(\mathcal{H}) \geq 2$):

i). $f = h(\|\cdot\|)$ for some increasing function $h : \mathbb{R}_+ \to \mathbb{R} \cup \{\infty\}$;

ii). For all perpendicular $x \perp y$, $f(x+y) \geq f(y)$;

iii). For all $z \in \mathcal{H}$, $P_f(z) = \lambda_z \cdot z$ for some $\lambda_z \in [0,1]$;

iv). $0 \in \operatorname{dom} f$ and $P_{f+\kappa} = P_f \circ P_\kappa$ for *all* positive homogeneous $\kappa$.

## i) $\implies$ iv)

Tree-structured group norms
(Jenatton et al.'11)

$$P_{\sum_i \|\cdot\|_{g_i}} = P_{\|\cdot\|_{g_1}} \circ \cdots \circ P_{\|\cdot\|_{g_k}}.$$

# Permutation Invariant $\Leftrightarrow$ Choquet Integral

$$\partial g(P_f(\mathbf{y})) \supseteq \partial g(\mathbf{y})$$

## Theorem

Let $f$ be permutation invariant and $g$ be the Choquet integral of some submodular set function $\mu$. Then, $P_{f-g} = P_f \circ P_g$.

# Permutation Invariant $\Leftrightarrow$ Choquet Integral

$$\partial g(\mathsf{P}_f(y)) \supseteq \partial g(y)$$

- $f$ permutation invariant $\Rightarrow \mathsf{P}_f(\mathbf{y}) /\!/ \mathbf{y}$:

$$(y_i - y_j)([\mathsf{P}_f(\mathbf{y})]_i - [\mathsf{P}_f(\mathbf{y})]_j) \geq 0$$

- $\partial g$ invariant to comonotone vectors
- Choquet integral (a.k.a. Lovász extension) of $\mu \quad 2^{[d]} - \mathbb{R}$:

$$g(\mathbf{w}) := \int_0^x \mu(\llbracket \mathbf{w} \geq t \rrbracket)\, dt - \int_{-\infty}^0 [\mu(\llbracket \mathbf{w} \geq t \rrbracket) - \mu(\llbracket d \rrbracket)]\, dt$$

## Theorem

Let $f$ be permutation invariant and $g$ be the Choquet integral of some submodular set function $\mu$. Then, $\mathsf{P}_{f+g} = \mathsf{P}_f \circ \mathsf{P}_g$.

# Permutation Invariant $\Leftrightarrow$ Choquet Integral

$$\partial g(\mathsf{P}_f(\mathbf{y})) \supseteq \partial g(\mathbf{y})$$

- $f$ permutation invariant $\Rightarrow \mathsf{P}_f(\mathbf{y}) /\!/ \mathbf{y}$:

$$(y_i - y_j)([\mathsf{P}_f(\mathbf{y})]_i - [\mathsf{P}_f(\mathbf{y})]_j) \geq 0$$

- $\partial g$ invariant to comonotone vectors
- Choquet integral (*a.k.a.* Lovász extension) of $\mu : 2^{[d]} \to \mathbb{R}$:

$$g(\mathbf{w}) := \int_0^\infty \mu(\llbracket \mathbf{w} \geq t \rrbracket)\, \mathrm{d}t + \int_{-\infty}^0 [\mu(\llbracket \mathbf{w} \geq t \rrbracket) - \mu([d])]\, \mathrm{d}t$$

## Theorem

Let $f$ be permutation invariant and $g$ be the Choquet integral of some submodular set function $\mu$. Then, $\mathsf{P}_{f-g} = \mathsf{P}_f \circ \mathsf{P}_g$.

# Permutation Invariant ⇔ Choquet Integral

$$\partial g(\mathsf{P}_f(\mathbf{y})) \supseteq \partial g(\mathbf{y})$$

- $f$ permutation invariant $\Rightarrow \mathsf{P}_f(\mathbf{y}) /\!/ \mathbf{y}$:

$$(y_i - y_j)([\mathsf{P}_f(\mathbf{y})]_i - [\mathsf{P}_f(\mathbf{y})]_j) \geq 0$$

- $\partial g$ invariant to comonotone vectors
- Choquet integral (a.k.a. Lovász extension) of $\mu : 2^{[d]} \to \mathbb{R}$:

$$g(\mathbf{w}) := \int_0^\infty \mu(\llbracket \mathbf{w} \geq t \rrbracket)\, \mathrm{d}t + \int_{-\infty}^0 [\mu(\llbracket \mathbf{w} \geq t \rrbracket) - \mu([d])]\, \mathrm{d}t$$

## Theorem

*Let $f$ be permutation invariant and $g$ be the Choquet integral of some submodular set function $\mu$. Then, $\mathsf{P}_{f+g} = \mathsf{P}_f \circ \mathsf{P}_g$.*
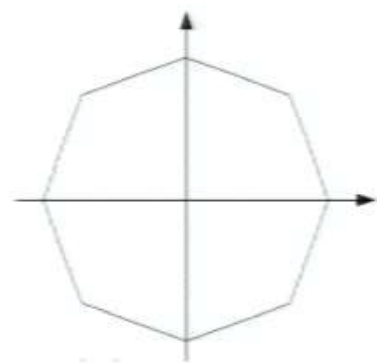
# Some Implications

> ## Theorem
>
> Let $f$ be permutation invariant and $g$ be the Choquet integral of some submodular set function. Then, $P_{f+g} = P_f \circ P_g$.

- Special case $f = \|\cdot\|_1$ in (Bach'11);
- $P_{\|\cdot\|_1 + \|\cdot\|_{TV}} = P_{\|\cdot\|_1} \circ P_{\|\cdot\|_{TV}}$ (Friedman et al.'07);
- $P_{\sum_{i=1}^k \|\cdot\|_{g_i}} = P_{\|\cdot\|_{g_1}} \circ \cdots \circ P_{\|\cdot\|_{g_k}}$ (Jenatton et al.'11)

# Some Implications

$$\|\mathbf{w}\|_{\mathrm{oscar}} = \sum_{i<j} \max\{|w_i|, |w_j|\}.$$

- Feature grouping (Bondell & Reich'08)
- $P_{\|\cdot\|_{\mathrm{oscar}}}$ in (Zhong & Kwok'11)

Let

$$\kappa_i(\mathbf{w}) := \sum_{j:j<i} \max\{|w_i|, |w_j|\}.$$

- $\|\mathbf{w}\|_{\mathrm{oscar}} = \sum_{i=2}^{d} \kappa_i(\mathbf{w})$
- $P_{\|\cdot\|_{\mathrm{oscar}}} = P_{\kappa_d} \circ \cdots \circ P_{\kappa_2}$
- Given $P_{\kappa_i}$, constant time for $P_{\kappa_{i+1}}$.

# Some Implications

> **Theorem**
>
> Let $f$ be permutation invariant and $g$ be the Choquet integral of some submodular set function. Then, $P_{f+g} = P_f \circ P_g$.

- Special case $f = \|\cdot\|_1$ in (Bach'11);
- $P_{\|\cdot\|_1 + \|\cdot\|_{TV}} = P_{\|\cdot\|_1} \circ P_{\|\cdot\|_{TV}}$ (Friedman et al.'07);
- $P_{\sum_{i=1}^{k} \|\cdot\|_{g_i}} = P_{\|\cdot\|_{g_1}} \circ \cdots \circ P_{\|\cdot\|_{g_k}}$ (Jenatton et al.'11)

# Some Implications

$$\|\mathbf{w}\|_{\mathrm{oscar}} = \sum_{i<j} \max\{|w_i|, |w_j|\}.$$

- Feature grouping (Bondell & Reich'08)
- $P_{\|\cdot\|_{\mathrm{oscar}}}$ in (Zhong & Kwok'11)

Let

$$\kappa_i(\mathbf{w}) := \sum_{j:j<i} \max\{|w_i|, |w_j|\}.$$

- $\|\mathbf{w}\|_{\mathrm{oscar}} = \sum_{i=2}^{d} \kappa_i(\mathbf{w})$
- $P_{\|\cdot\|_{\mathrm{oscar}}} = P_{\kappa_d} \circ \cdots \circ P_{\kappa_2}$
- Given $P_{\kappa_i}$, constant time for $P_{\kappa_{i+1}}$.

# Sufficient Condition Fails?

- (Martins et al.'11) showed that, under a shrinkage assumption, the prox-decomposition (even not true) can still be used in a subgradient-type algorithm

- (Yu'13a) showed that a simple linearization of the proximal map, i.e.

$$P_{\sum_k f_k} \approx \sum_k P_{f_k},$$

yields slightly faster convergence than the smoothing trick

# Summary

- Posed the question: $P_{f+g} \overset{?}{=} P_f \circ P_g \overset{?}{=} P_g \circ P_f$;
- Presented a sufficient condition: $\partial g(P_f(\mathbf{y})) \supseteq \partial g(\mathbf{y})$;
- Identified two major cases;
- Immediately useful if plugged into PG;

# Summary

- Posed the question: $P_{f+g} \overset{?}{=} P_f \circ P_g \overset{?}{=} P_g \circ P_f$;
- Presented a sufficient condition: $\partial g(P_f(\mathbf{y})) \supseteq \partial g(\mathbf{y})$;
- Identified two major cases;
- Immediately useful if plugged into PG;

# Thanks!

# Non-Uniform Camera Shake Removal using a Spatially Adaptive Sparse Penalty

Haichao Zhang [1,2]      David Wipf [3]

NIPS, December 2013

1

# Problem & Objective

- Problem
  - **Camera shake blur** caused by **relative movement** between camera and scene **during exposure**.

# Problem & **Objective**

- Objective
  - Recover the **sharp image** from a single **blurry** image with **unknown** camera shake.

# Challenge I
## Ill-posed Problem: no unique solution

Sharp image $\mathbf{x}$  Blur kernel $\mathbf{h}$

$$y = h * x + n$$

=   * 

**No Blur Solution**

**Blurry image**

=   * 

=   * 

**True Solution**

## Observation Model [Projective Motion Path]

$$y = \sum_j w_j P_j x + n = Dw + n \qquad D = [P_1 x, P_2 x, \cdots]$$

blurry image · blur vector · sharp image · noise



$$= [ \qquad \cdots ] w$$

Tai et al. 2009
Hirsch et al. ICCV 2011

$$D = [P_1 x, P_2 x, \cdots]$$

6

# Non-Uniform Observation Model
## ... regarding challenge II

**Observation Model** [Projective Motion Path]

$$y = \sum_j w_j P_j x + n = Dw + n \qquad D = [P_1 x, P_2 x, \cdots]$$

$$= Hx + n \qquad H = \sum_j w_j P_j$$

blurry image   blur vector   sharp image   noise

$$= [\phantom{P_1x,P_2x}\cdots] \, w = [\phantom{h_1,h_2}\cdots]$$

Tai et al. 2009
Hirsch et al. ICCV 2011

$$D = [P_1 x, P_2 x, \cdots] \qquad H = [h_1, h_2, \cdots]$$

7

**Likelihood**

$$p(\mathbf{y}|\mathbf{x}, \mathbf{w}) \propto \exp\left\{-\frac{1}{\lambda}\|\mathbf{y} - \mathbf{H}\mathbf{x}\|_2^2\right\}$$

**Image Prior [sparse gradient]**

$$p(\mathbf{x}) \propto \exp\left[-\sum_i g(x_i)\right]$$

$g(x_i)$ is a concave function

**Blur Prior**     $p(\mathbf{w})$



Work in derivative domain

$\mathbf{x}$   (vectorized) derivatives of the sharp image

$\mathbf{y}$   (vectorized) derivatives of the blurry image

Fergus et al. Removing camera shake from a single image, SIGGRAPH'06

8

# Direct MAP ?

## MAP Estimation

$$\max_{\mathbf{x},\mathbf{w}\geq 0} p(\mathbf{x},\mathbf{w}|\mathbf{y}) \equiv \min_{\mathbf{x},\mathbf{w}\geq 0} -\log[p(\mathbf{y}|\mathbf{x},\mathbf{w})p(\mathbf{x})p(\mathbf{w})]$$

$$\min_{\mathbf{x},\mathbf{w}\geq 0} \|\mathbf{y}-\mathbf{Hx}\|_2^2 + \alpha \sum_i g(x_i) + \beta \sum_j f(w_j)$$

– Local minima and "no-blur" solution
– Empirical tricks
  • initialization, structure selection/prediction  [Cho and Lee SIGGRAPH Asia'09, Xu & Jia ECCV'10, Hu et al. BMVC'12]

# Type-II Estimation

**Likelihood**     $p(\mathbf{y}|\mathbf{x}, \mathbf{w}) \propto \exp\left\{-\frac{1}{\lambda}\|\mathbf{y} - \mathbf{H}\mathbf{x}\|_2^2\right\}$

**Image Prior**     $p(\mathbf{x}) \sim \mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{\Gamma})$     $\mathbf{\Gamma} \triangleq \text{diag}[\gamma]$

$$\max_{\gamma, \mathbf{w}, \lambda \geq 0} \int_{\mathbf{x}} p(\mathbf{y}|\mathbf{x}, \mathbf{w}) p(\mathbf{x}) d\mathbf{x} \qquad \text{uniform } p(\mathbf{w}) \text{ is used}$$

Tipping et al. Sparse Bayesian Learning and Relevance Vector Machine, JMLR 2011
Levin et al. Efficient Marginal Likelihood Optimization in Blind Deconvolution, CVPR 2011
Wipf et al., Latent Variable Bayesian Model for Promoting Sparsity, IEEE TIT, 2011

# Effective Cost Function

## The Cost Function

$$\min_{\mathbf{x};\gamma,\mathbf{w},\lambda\geq 0} \frac{1}{\lambda}\|\mathbf{y}-\mathbf{Hx}\|_2^2 + \sum_i \psi(|x_i|\|\mathbf{h}_i\|_2, \lambda) + (n-m)\log\lambda$$

$$\psi(u,\lambda) \triangleq \frac{2u}{u+\sqrt{4\lambda+u^2}} + \log\left(2\lambda + u^2 + u\sqrt{4\lambda+u^2}\right) \quad u \geq 0$$

Wipf et al., Latent Variable Bayesian Model for Promoting Sparsity, IEEE TIT, 2011

# Effective Cost Function

## The Cost Function

$$\min_{\mathbf{x};\gamma,\mathbf{w},\lambda\geq 0} \boxed{\frac{1}{\lambda}\|\mathbf{y}-\mathbf{H}\mathbf{x}\|_2^2} + \sum_i \psi(|x_i|\|\mathbf{h}_i\|_2, \lambda) + (n-m)\log\lambda$$

$$\psi(u,\lambda) \triangleq \frac{2u}{u+\sqrt{4\lambda+u^2}} + \log\left(2\lambda+u^2+u\sqrt{4\lambda+u^2}\right) \quad u \geq 0$$

**reconstruction error**

Wipf et al., Latent Variable Bayesian Model for Promoting Sparsity, IEEE TIT, 2011

# Effective Cost Function

## The Cost Function

$$\min_{\mathbf{x};\gamma,\mathbf{w},\lambda\geq0} \frac{1}{\lambda}\|\mathbf{y}-\mathbf{Hx}\|_2^2 + \boxed{\sum_i \psi(|x_i|\|\mathbf{h}_i\|_2,\lambda)} + (n-m)\log\lambda$$

$$\psi(u,\lambda) \triangleq \frac{2u}{u+\sqrt{4\lambda+u^2}} + \log\left(2\lambda+u^2+u\sqrt{4\lambda+u^2}\right) \quad u\geq0$$

**sparse penalty function**

$\psi(u)$ is a concave, non-decreasing function of $u$

Wipf et al., Latent Variable Bayesian Model for Promoting Sparsity, IEEE TIT, 2011

# Effective Cost Function

## The Cost Function

$$\min_{\mathbf{x};\gamma,\mathbf{w},\lambda\geq 0} \frac{1}{\lambda}\|\mathbf{y} - \mathbf{Hx}\|_2^2 + \sum_i \psi(|x_i|\|\mathbf{h}_i\|_2, \lambda) + \boxed{(n-m)\log\lambda}$$

$$\psi(u,\lambda) \triangleq \frac{2u}{u + \sqrt{4\lambda + u^2}} + \log\left(2\lambda + u^2 + u\sqrt{4\lambda + u^2}\right) \quad u \geq 0$$

**noise level penalty term**

Wipf et al., Latent Variable Bayesian Model for Promoting Sparsity, IEEE TIT, 2011

# Effective Cost Function

## The Cost Function

$$\min_{\mathbf{x};\gamma,\mathbf{w},\lambda\geq 0} \frac{1}{\lambda}\|\mathbf{y}-\mathbf{Hx}\|_2^2 + \sum_i \psi(|x_i|\|\mathbf{h}_i\|_2, \lambda) + (n-m)\log\lambda$$

$$\psi(u,\lambda) \triangleq \frac{2u}{u+\sqrt{4\lambda+u^2}} + \log\left(2\lambda+u^2+u\sqrt{4\lambda+u^2}\right) \quad u \geq 0$$

can be solved using the *majorization-minimization* technique

Wipf et al., Latent Variable Bayesian Model for Promoting Sparsity, IEEE TIT, 2011

# Effective Cost Function

## The Cost Function

$$\min_{\mathbf{x};\gamma,\mathbf{w},\lambda\geq 0} \frac{1}{\lambda}\|\mathbf{y}-\mathbf{Hx}\|_2^2 + \sum_i \psi(|x_i|\|\mathbf{h}_i\|_2, \lambda) + (n-m)\log\lambda$$

$$\psi(u,\lambda) \triangleq \frac{2u}{\cdots} + \log\left(2\lambda + u^2 + u\sqrt{4\lambda + u^2}\right) \quad u \geq 0$$

**looks similar, what's the real advantage ... (over the regular MAP)?**

can be solved using the majorization-minimization technique

$$\min_{\mathbf{x},\mathbf{w}\geq 0} \frac{1}{\lambda}\|\mathbf{y}-\mathbf{Hx}\|_2^2 + \alpha\sum_i g(x_i) + \beta\sum_j f(w_j)$$

- **Effect of Spatially-Variant Blur on H**
  - Imbalanced Column of **H** ( $H = [h_1, h_2, \cdots]$ )
    - Each column of **H** corresponds to a localized blur kernel
    - Large blur has smaller L2 norm ($h_i \geq 0, \sum h_i = 1$)
    - Columns of **H** have different L2 norms (local kernel norm $\|h_i\|_2$)

$$H = [h_1, h_2, \cdots]$$

- **Effect of Spatially-Variant Blur on H**
  - Imbalanced Column of **H** ( $H = [h_1, h_2, \cdots]$ )
    - Each column of **H** corresponds to a localized blur kernel
    - Large blur has smaller L2 norm ($h_i \geq 0, \sum h_i = 1$)
    - Columns of **H** have different L2 norms (local kernel norm $\|h_i\|_2$)

$$H = [h_1, h_2, \cdots]$$



**Bias image recovery and therefore affect the kernel estimation.**

- **Column-Normalized Sparse Estimation**

$$\min_{\mathbf{x};\gamma,\mathbf{w},\lambda\geq0} \frac{1}{\lambda}\|\mathbf{y}-\mathbf{Hx}\|_2^2 + \sum_i \psi(|x_i|\|\mathbf{h}_i\|_2, \lambda) + (n-m)\log\lambda$$

**local kernel norm embedded**

**compensates for the spatial variance**

- **Column-Normalized Sparse Estimation**

$$\min_{\mathbf{x};\gamma,\mathbf{w},\lambda\geq0} \frac{1}{\lambda}\|\mathbf{y}-\mathbf{H}\mathbf{x}\|_2^2 + \sum_i \psi(\boxed{|x_i|\|\mathbf{h}_i\|_2},\lambda) + (n-m)\log\lambda$$

$$z_i \triangleq x_i\|\mathbf{h}_i\|_2$$

$$\min_{\mathbf{z};\gamma,\mathbf{w},\lambda\geq0} \frac{1}{\lambda}\|\mathbf{y}-\tilde{\mathbf{H}}\mathbf{z}\|_2^2 + \sum_i \psi(|z_i|,\lambda) + (n-m)\log\lambda$$

$\tilde{\mathbf{H}}$ is the column-normalized $\mathbf{H}$

# Model Properties
## Automated Column-Normalization

- **Column-Normalized Sparse Estimation**

$$\min_{\mathbf{x};\gamma,\mathbf{w},\lambda\geq0} \frac{1}{\lambda}\|\mathbf{y}-\mathbf{Hx}\|_2^2 + \sum_i \psi(\boxed{|x_i|\|\mathbf{h}_i\|_2}, \lambda) + (n-m)\log\lambda$$

$$z_i \triangleq x_i\|\mathbf{h}_i\|_2$$

*large structure, low blur region wil be naturally emphasized*

$$\min_{\mathbf{z};\gamma,\mathbf{w},\lambda\geq0} \frac{1}{\lambda}\|\mathbf{y}-\tilde{\mathbf{H}}\mathbf{z}\|_2^2 + \sum_i \psi(|z_i|, \lambda) + (n-m)\log\lambda$$

**Avoids premuture favoring of any one element of z over another
(thus avoid biased image recovery )**

# Model Properties
## Automated Column-Normalization

- **Effects of Column-Normalization** (blind deblurring)



Blurry image from Harmeling et al.
*NIPS* 2010

- **Two Effects of Blur on Sparsity Measure** (Lp-norm)

  1. Reduces sparsity $\qquad \sum_i |y_i|^p \nearrow$

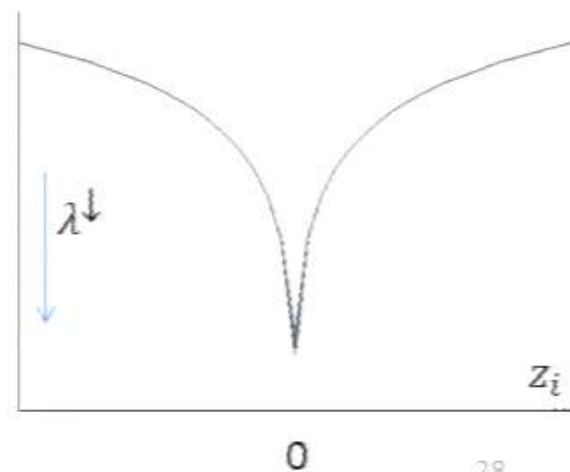  2. Reduces variance $\qquad \sum_i |y_i|^p \searrow$



signal

original
blurred

Levin et al., *CVPR* 2009



sparse penalty val.

original
blurred

factor 2
dominates

natural image prior

$\sum_i |y_i|^p$

$p$

25

- **Two Effects of Blur on Sparsity Measure** (Lp-norm)

  1. Reduces sparsity $\quad \sum_i |y_i|^p \nearrow$

  2. Reduces variance $\quad \sum_i |y_i|^p \searrow$

**Very concave penalty (e.g. L0-norm) should be used**

sparse penalty val.

factor 1 dominates

factor 2 dominates

natural image prior

original
blurred

signal

original
blurred

$\sum |y_i|^p$

$p$

Levin et al., *CVPR 2009*

26

- **Two Effects of Blur on Sparsity Measure** (Lp-norm)

1. Reduces sparsity $\quad \sum_i |y_i|^p \nearrow$

2. Reduces variance $\quad \sum_i |y_i|^p \searrow$

**Very concave penalty (e.g. L0-norm) should be used**

**Non-Convex Problem!!!**



sparse penalty val.

factor 1 dominates

factor 2 dominates

natural image prior

original
blurred

$\sum |y_i|^p$

$p$

Levin et al., *CVPR* 2009

- **The penalty function in the proposed model**
  - A qualified "very concave" sparse penalty

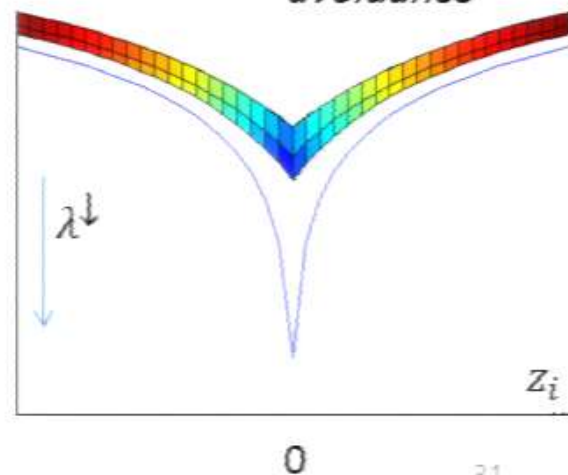$$\text{As } \lambda \rightarrow 0 \, , \; \sum \psi(|z_i|, \lambda) \rightarrow C\|z\|_0$$

*no-blur solution avoidance*

$$\psi(u, \lambda) \triangleq \frac{2u}{u + \sqrt{4\lambda + u^2}} + \log\left(2\lambda + u^2 + u\sqrt{4\lambda + u^2}\right) \quad u \geq 0$$

- **The penalty function in the proposed model**
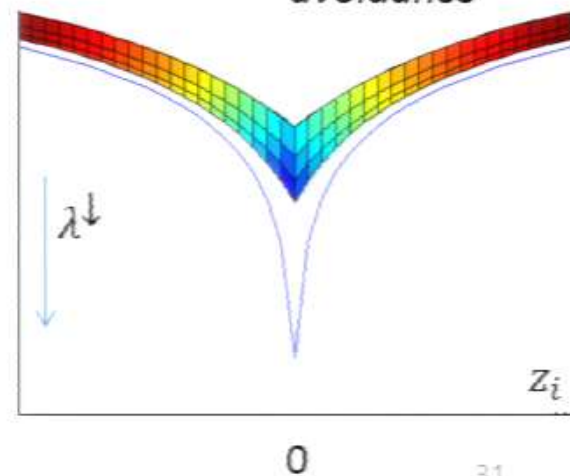  - A qualified "very concave" sparse penalty

    As $\lambda \longrightarrow 0$, $\sum \psi(|z_i|, \lambda) \longrightarrow C\|z\|_0$    *no-blur solution avoidance*

  - Adaptive penalty shape

    As $\lambda$ is large, $\sum \psi(|z_i|, \lambda) \longrightarrow 2\|z\|_1/\sqrt{\lambda}$

$$\psi(u, \lambda) \triangleq \frac{2u}{u + \sqrt{4\lambda + u^2}} + \log\left(2\lambda + u^2 + u\sqrt{4\lambda + u^2}\right) \quad u \geq 0$$

- **The penalty function in the proposed model**
  - A qualified "very concave" sparse penalty

    $$\text{As } \lambda \longrightarrow 0 \text{ , } \sum \psi(|z_i|, \lambda) \longrightarrow C\|z\|_0$$

    *no-blur solution avoidance*

  - Adaptive penalty shape

    $$\text{As } \lambda \text{ is large, } \sum \psi(|z_i|, \lambda) \longrightarrow 2\|z\|_1/\sqrt{\lambda}$$

    *local-minia avoidance*

    $$\text{If } \lambda_1 < \lambda_2, \text{ then } \psi(u, \lambda_1) < \psi(u, \lambda_2) \\ \text{for } u \geq 0$$
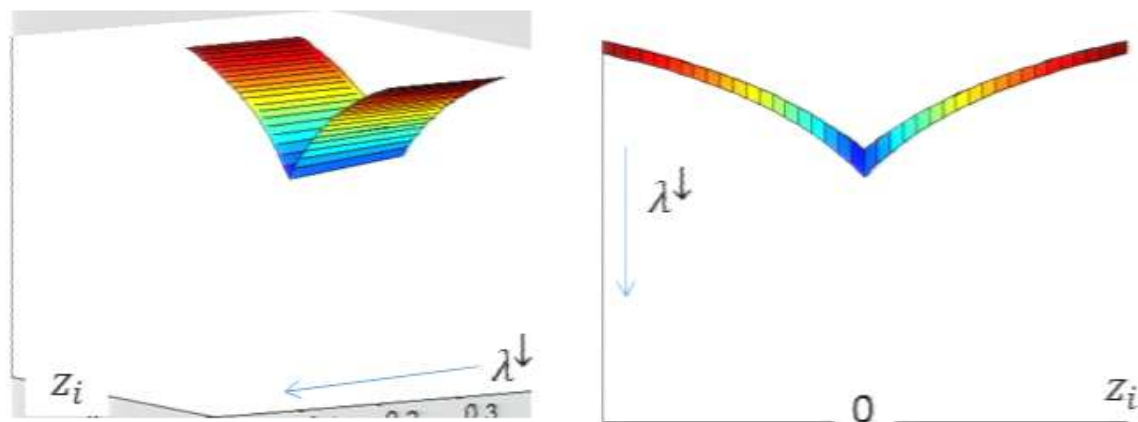
$$\psi(u, \lambda) \triangleq \frac{2u}{u + \sqrt{4\lambda + u^2}} + \log\left(2\lambda + u^2 + u\sqrt{4\lambda + u^2}\right) \quad u \geq 0$$

- **The penalty function in the proposed model**
  - A qualified "very concave" sparse penalty

    As $\lambda \longrightarrow 0$ , $\sum \psi(|z_i|, \lambda) \longrightarrow C\|z\|_0$     *no-blur solution avoidance*

  - Adaptive penalty shape

    As $\lambda$ *is large,* $\sum \psi(|z_i|, \lambda) \longrightarrow 2\|z\|_1/\sqrt{\lambda}$     *local-minia avoidance*

    If $\lambda_1 < \lambda_2$, then $\psi(u, \lambda_1) < \psi(u, \lambda_2)$
    for $u \geq 0$

$$\psi(u, \lambda) \triangleq \frac{2u}{u + \sqrt{4\lambda + u^2}} + \log\left(2\lambda + u^2 + u\sqrt{4\lambda + u^2}\right) \quad u \geq 0$$



Wipf et al., Latent Varianbel Bayesian Model for Promoting Sparsity, IEEE TIT, 2011

31

- **The penalty function in the proposed model**
  - A qualified "very concave" sparse penalty

    *no-blur solution avoidance*

    $$\text{As } \lambda \rightarrow 0, \ \sum \psi(|z_i|, \lambda) \rightarrow C\|z\|_0$$

  - Adaptive penalty shape

    $$\text{As } \lambda \text{ is large, } \sum \psi(|z_i|, \lambda) \rightarrow 2\|z\|_1/\sqrt{\lambda}$$

    *local-minia avoidance*

    $$\text{If } \lambda_1 < \lambda_2, \text{ then } \psi(u, \lambda_1) < \psi(u, \lambda_2)$$
    $$\text{for } u \geq 0$$



$$\psi(u, \lambda) \triangleq \frac{2u}{u + \sqrt{4\lambda + u^2}} + \log\left(2\lambda + u^2 + u\sqrt{4\lambda + u^2}\right) \quad u \geq 0$$

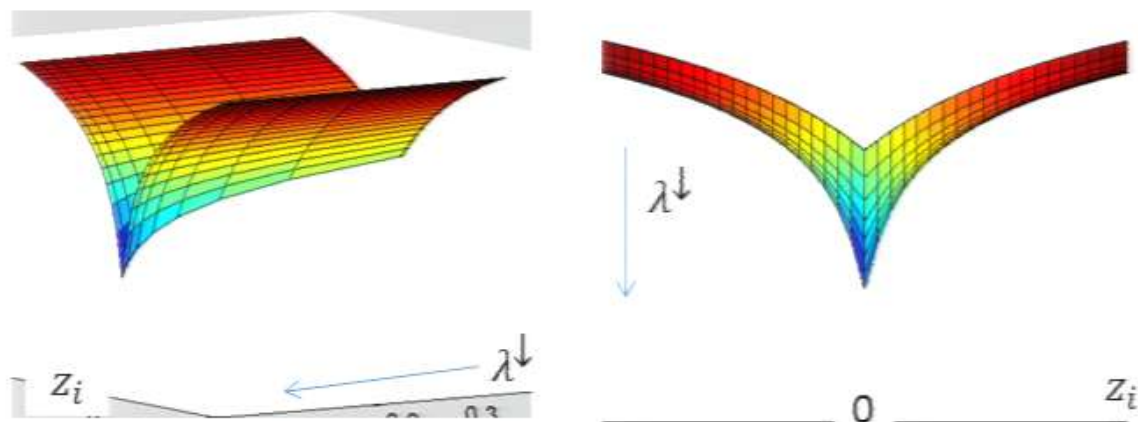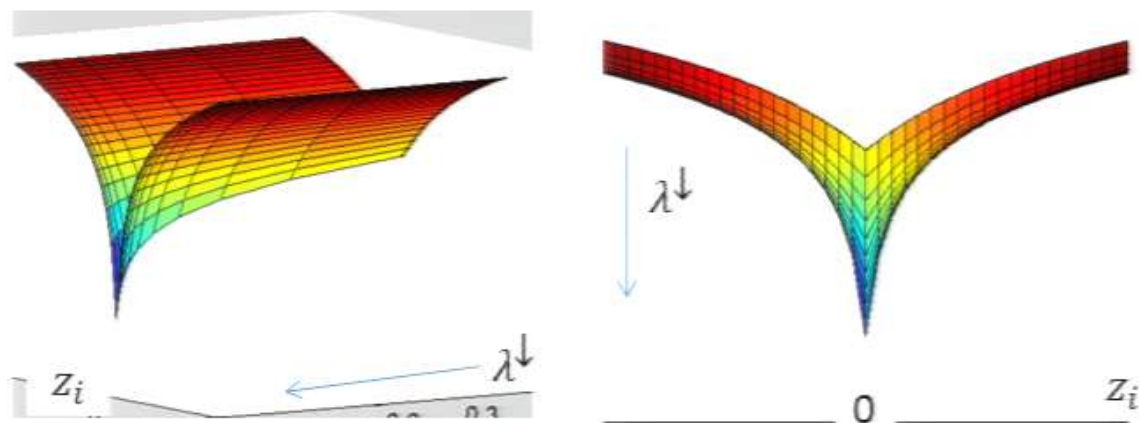Wipf et al., Latent Varianbel Bayesian Model for Promoting Sparsity, IEEE TIT, 2011

- **Implications on Camera Shake Removal**
  - Initially, $\lambda$ is large, penalty function is less concave
    - de-emphasize high blur regions ($z_i$ small)     $z_i = x_i \|h_i\|_2$
    - focus first on large structure ($x_i$ large), low blur ($\|h_i\|_2$ large) regions
  - Later, $\lambda$ is reduced, relative concavity of $\psi$ is increased, more fine details will be recovered
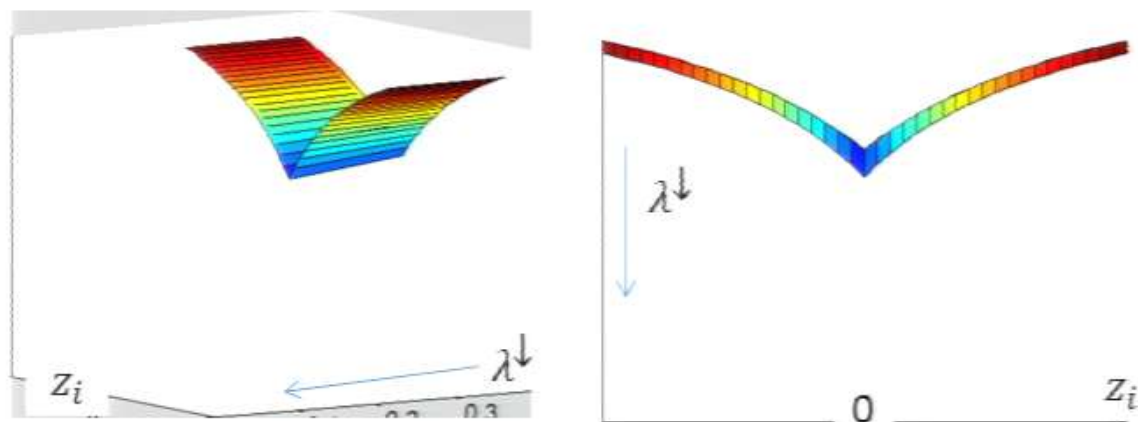


32

- **Implications on Camera Shake Removal**
  - Initially, $\lambda$ is large, penalty function is less concave
    - de-emphasize high blur regions ($z_i$ small) $\quad z_i = x_i \|h_i\|_2$
    - focus first on large structure ($x_i$ large), low blur ($\|h_i\|_2$ large) regions
  - Later, $\lambda$ is reduced, relative concavity of $\psi$ is increased, more fine details will be recovered



32

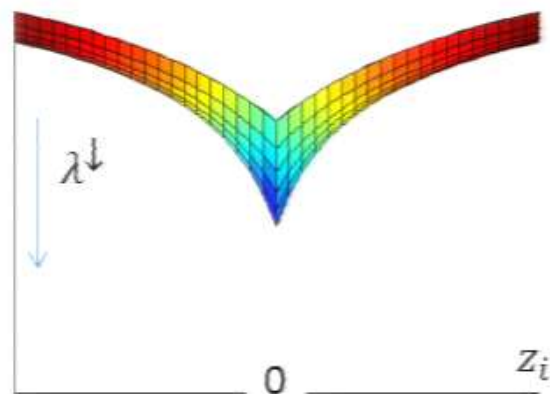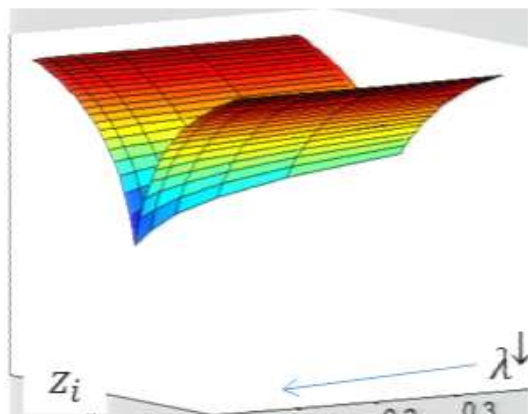# Model Properties
## Noise Dependent Homotopy Continuation

- **Implications on Camera Shake Removal**
  - Initially, $\lambda$ is large, penalty function is less concave
    - de-emphasize high blur regions ($z_i$ small)       $z_i = x_i \| h_i \|_2$
    - focus first on large structure ($x_i$ large), low blur ($\| h_i \|_2$ large) regions
  - Later, $\lambda$ is reduced, relative concavity of $\psi$ is increased, more fine details will be recovered

- **Implications on Camera Shake Removal**
  - Initially, $\lambda$ is large, penalty function is less concave
    - de-emphasize high blur regions ($z_i$ small)  $\quad z_i = x_i \| h_i \|_2$
    - focus first on large structure ($x_i$ large), low blur ($\| h_i \|_2$ large) regions
  - Later, $\lambda$ is reduced, relative concavity of $\psi$ is increased, more fine details will be recovered
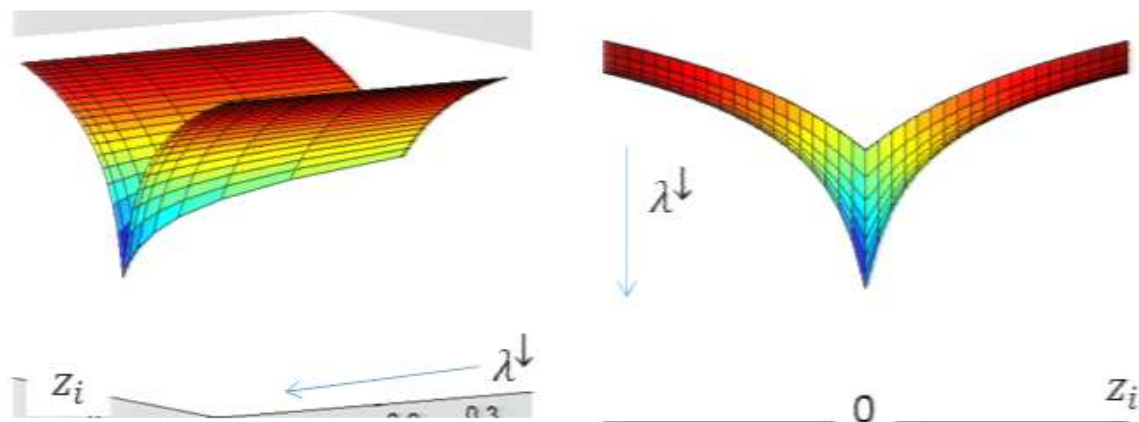
- **Implications on Camera Shake Removal**

  - Initially, $\lambda$ is large, penalty function is less concave
    - de-emphasize high blur regions ($z_i$ small)     $z_i = x_i \|h_i\|_2$
    - focus first on large structure ($x_i$ large), low blur ($\|h_i\|_2$ large) regions

  - Later, $\lambda$ is reduced, relative concavity of $\psi$ is increased, more fine details will be recovered
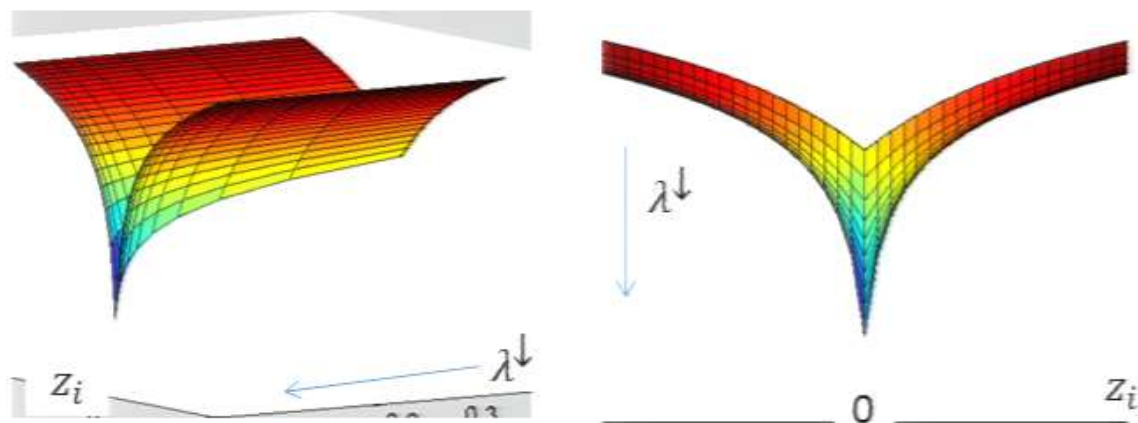
- **Implications on Camera Shake Removal**
  - Initially, $\lambda$ is large, penalty function is less concave
    - de-emphasize high blur regions ($z_i$ small)     $z_i = x_i \|h_i\|_2$
    - focus first on large structure ($x_i$ large), low blur ($\|h_i\|_2$ large) regions
  - Later, $\lambda$ is reduced, relative concavity of $\psi$ is increased, more fine details will be recovered

- **Implications on Camera Shake Removal**

  - Initially, $\lambda$ is large, penalty function is less concave

    - de-emphasize high blur regions ($z_i$ small)   $z_i = x_i \|h_i\|_2$

    - focus first on large structure ($x_i$ large), low blur ($\|h_i\|_2$ large) regions

  - Later, $\lambda$ is reduced, relative concavity of $\psi$ is increased, more fine details will be recovered

- **Implications on Camera Shake Removal**
  - Initially, $\lambda$ is large, penalty function is less concave
    - de-emphasize high blur regions ($z_i$ small)     $z_i = x_i \|h_i\|_2$
    - focus first on large structure ($x_i$ large), low blur ($\|h_i\|_2$ large) regions
  - Later, $\lambda$ is reduced, relative concavity of $\psi$ is increased, more fine details will be recovered



32

- ## **Implications on Camera Shake Removal**

  - Initially, $\lambda$ is large, penalty function is less concave

    - de-emphasize high blur regions ($z_i$ small)   $z_i = x_i \|h_i\|_2$

    - focus first on large structure ($x_i$ large), low blur ($\|h_i\|_2$ large) regions

  - Later, $\lambda$ is reduced, relative concavity of $\psi$ is increased, more fine details will be recovered



$\lambda \downarrow$

$z_i$   $\lambda \downarrow$

$z_i$

0   $z_i$

- **The proposed cost function**

$$\min_{\mathbf{z};\gamma,\mathbf{w},\lambda \geq 0} \frac{1}{\lambda}\|\mathbf{y} - \tilde{\mathbf{H}}\mathbf{z}\|_2^2 + \sum_i \psi(|z_i|,\lambda) + (n-m)\log\lambda$$

  – Learning $\lambda$
  – Tuning parameter free

# Experiments

- Test Images
  - Real-world blurry images from literature
- Compared Methods
  - Harmeling et al. *NIPS* 2010
  - Whyte et al.      *CVPR* 2010
  - Gupta et al.      *ECCV* 2010
  - Hirsch et al.      *ICCV* 2011
  - Joshi et al.   *SIGGRAPH* 2010 [hardware asisted]
  - Cho et al.    *Pacific Graphics* 2012  [dual image]

*All the compared results are from the original authors*

# Experimental Results
## An illustration

A test blurry image from
Harmeling et al. , *NIPS* 2010.



Blurry Image

# Experimental Results
## An illustration



Estimated Kernel Patterns

Our

# Experimental Results
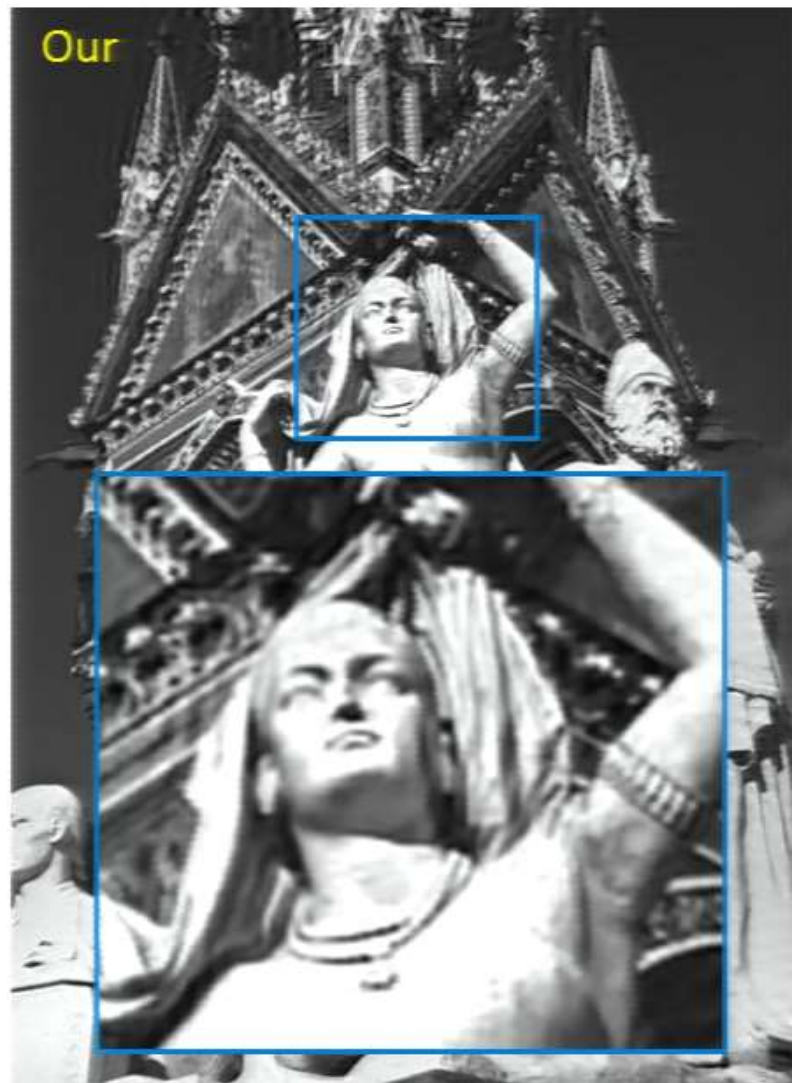## comparison with Harmeling *et al.* NIPS'10

# Experimental Results
## comparison with Harmeling *et al.* NIPS'10

Harmeling *et al.*

Our

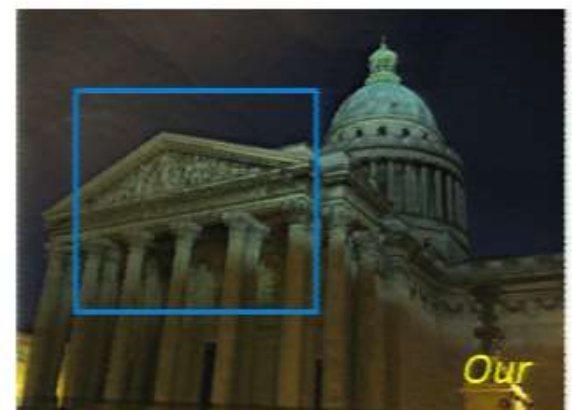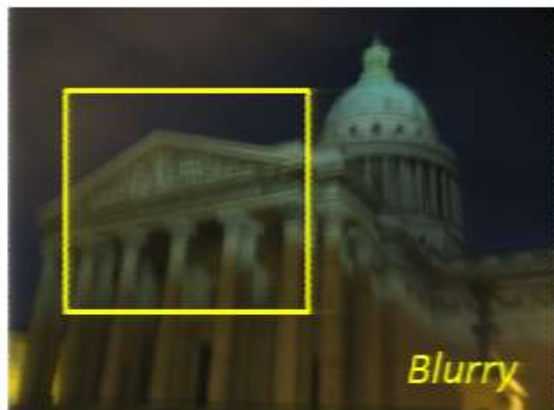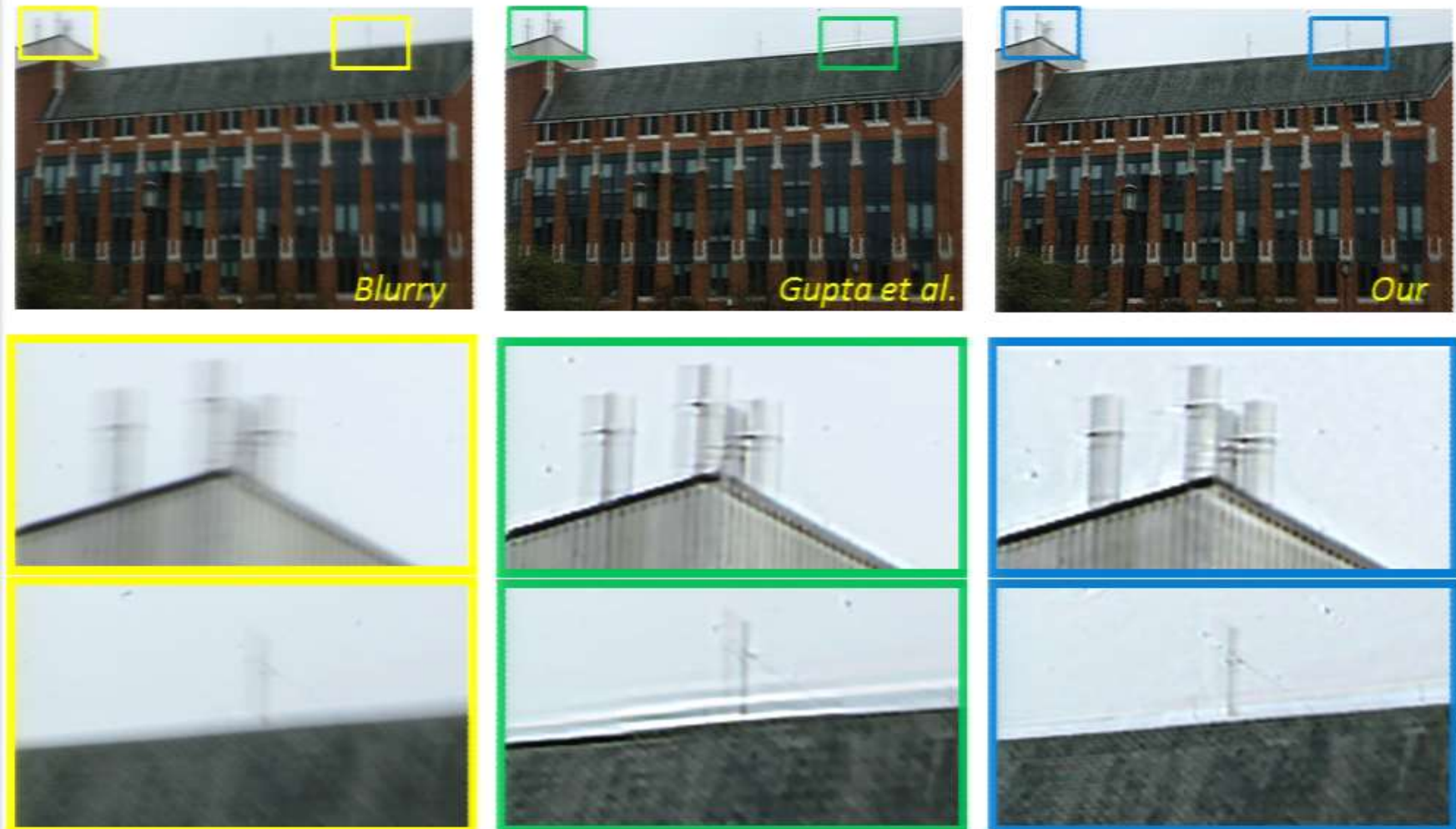# Experimental Results
## comparison with Whyte *et al.* CVPR'10



Blurry

Whyte et al.

Our

O. Whyte et al., *Non-uniform deblurring for shaken images*, CVPR, 2010.
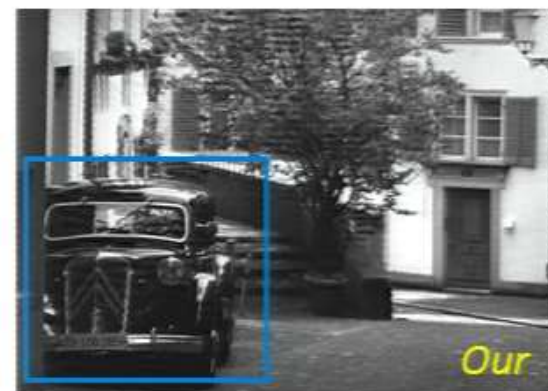
39

# Experimental Results
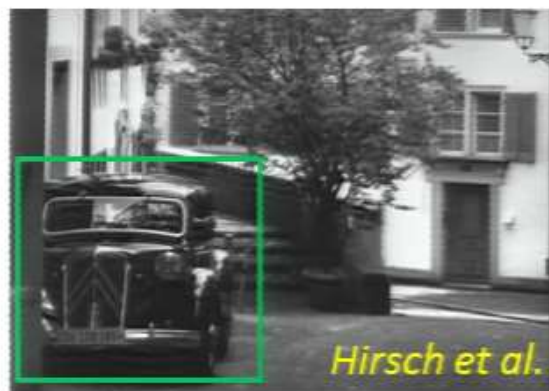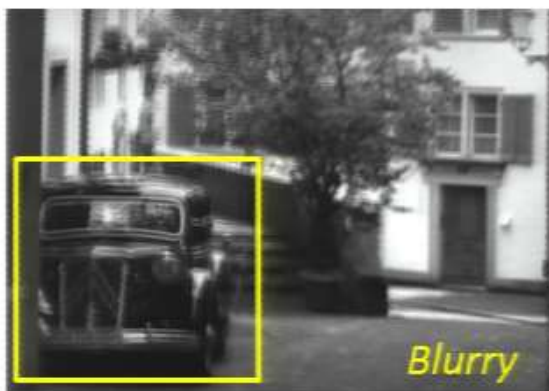## comparison with Gupta *et al.* ECCV'10

Gupta et al., *Single image deblurring using motion density functions*, ECCV, 2010.

40

# Experimental Results
## comparison with Hirsch *et al.* ICCV'11

Blurry

Hirsch et al.

Our

S. Hirsch et al. , *Fast removal of non-uniform camera shake*, ICCV, 2011.

# Experimental Results
## comparison with Joshi *et al.* SIGGRAPH'10



Blurry

Joshi et al.

Our

[hardware asisted]

N. Joshi et al. , *Image deblurring using inertial measurement sensors*, SIGGRAPH, 2010.

# Experimental Results
## comparison with Cho *et al.* PG'12

Cho et al., *Registration Based Non-uniform Motion Deblurring*, Pacific Graphics, 2012.

# Experimental Results
## comparison with Cho *et al.* PG'12

Blurry I

Blurry II

Life is gooood

Life is gooood

Cho    [image pair-based]

Our

Try something new

Try something new

get away from you
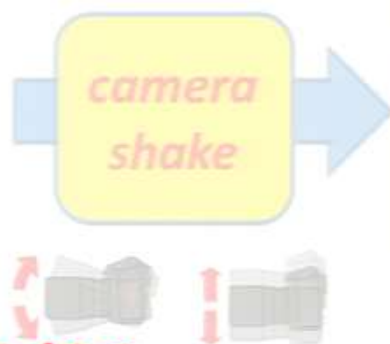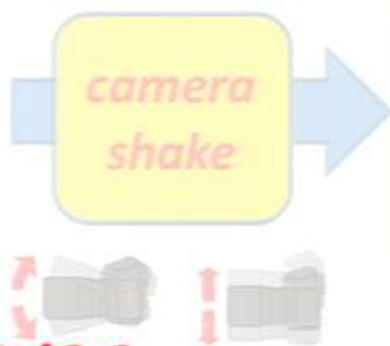
get away from you

# Summary

- An effective approach for camera shake removal
  - simple & clear cost function
- Model property analysis
  - automated column normalization (spatially adaptive sparsity): high-bur, low structure regions will be de-emphasized first, and emphasized progressively later
  - noise dependent homotopy continuation
  - tuning parameter free
- State-of-the-art performance on real-images
- Applicable to other problems (e.g., structured dictionary learning)

# Thank you!

# Questions?

*camera shake*

**Welcome to our poster Fri26**

Research

Experimental Results
comparison with Cho *et al.* PG'12

# Summary

- **An effective approach for camera shake removal**
  - simple & clear cost function
- **Model property analysis**
  - automated column normalization (spatially adaptive sparsity): high-bur, low structure regions will be de-emphasized first, and emphasized progressively later
  - noise dependent homotopy continuation
  - tuning parameter free
- **State-of-the-art performance on real-images**
- **Applicable to other problems (e.g., structured dictionary learning)**