

# Optimization

Computer Vision, Lecture 13  
Michael Felsberg  
Computer Vision Laboratory  
Department of Electrical Engineering



1

## Optimization: Overview

Function	Set	Output (codomain / target set)	
		Continuous	Discrete
Input (domain of definition)	Continuous	Lecture 15	Lecture 15
	Discrete	Lecture 13	Lecture 13

e.g.: stereo

e.g.: segmentation



2

## Why Optimization?

- Computer vision algorithms are usually very complex
  - Many parameters (dependent)
  - Data dependencies (non-linear)
  - Outliers and occlusions (noise)
- Classical approach
  - Trial and error (hackers' approach)
  - Encyclopedic knowledge (recipes)
  - Black-boxes + glue (hide problems)



3

## Why Optimization?

- Establishing CV as scientific discipline
  - Derive algorithms from first principles (*optimal solution*)
  - Automatic choice of parameters (*parameter free*)
  - Systematic evaluation (*benchmarks on standard datasets*)



4

5

## Optimization: howto

1. Choose a *scalar* measure (objective function) of success
  - From the benchmark
  - Such that optimization becomes *feasible*
  - Project functionality onto *one dimension*
2. Approximate the world with a model
  - Definition: allows to make *predictions*
  - Purpose: makes optimization *feasible*
  - Enables: *proper* choice of dataset

Similar to  
economics  
(money rules)



5

6

## Optimization: howto

3. Apply suitable framework for model fitting
  - This lecture
  - Systematic part (1 & 2 are ad hoc)
  - Current focus of research
4. Analyze resulting algorithm
  - Find *appropriate* dataset
  - Ignore runtime behavior (*highly non-optimized Matlab code* ;-)



6

7

## Examples

- Relative pose (F-matrix) estimation:
  - Algebraic error (quadratic form)
  - Linear solution by SVD
  - Robustness by random sampling (RANSAC)
  - Result: F and inlier set
- Bundle adjustment
  - Geometric (reprojection) error (quadratic error)
  - Iterative solution using LM
  - Result: camera pose and 3D points



7

8

## Taxonomy

- Objective function
  - Domain/manifold (algebraic error, geometric error, data dependent)
  - Robustness (explicitly in error norm, implicitly by Monte-Carlo approach)
- Model / simplification
  - Linearity (limited order), Markov property, regularization
- Algorithm
  - Approximate / analytic solutions (minimal problem)
  - Minimal solutions (over-determined)



8

9

## Taxonomy example: KLT

- Objective function
  - Domain/manifold: grey values / RGB / ...
  - Robustness: no (quadratic error, no regularization)

$$\varepsilon(\mathbf{d}) = \sum_{\mathbf{x} \in \mathcal{N}} w(\mathbf{x}) |f(\mathbf{x} - \mathbf{d}) - g(\mathbf{x})|^2$$

- Model: Brightness constancy, image shift

$$f(\mathbf{x} - \mathbf{d}) = g(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{N}$$

- local linearization (Taylor expansion)

$$f(\mathbf{x} - \mathbf{d}) \approx f(\mathbf{x}) - \mathbf{d}^T \nabla f(\mathbf{x}) \quad \nabla f = \left[ \frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \right]^T$$



9

10

## Taxonomy: KLT

- Algorithm
  - iterative solution of normal equations (Gauss-Newton)

$$\left( \sum_{\mathcal{R}} w(\mathbf{x}) \nabla f(\mathbf{x}) \nabla^T f(\mathbf{x}) \right) \mathbf{d} = \sum_{\mathcal{R}} w(\mathbf{x}) \nabla f(\mathbf{x}) (f(\mathbf{x}) - g(\mathbf{x}))$$

$$\mathbf{T} \mathbf{d} = \mathbf{r}$$

- $\mathbf{T}$ : structure tensor (orientation tensor from outer product of gradients)

$$\nabla f \nabla^T f = \begin{bmatrix} \left( \frac{\partial f}{\partial x} \right)^2 & \frac{\partial f}{\partial x} \frac{\partial f}{\partial y} \\ \frac{\partial f}{\partial x} \frac{\partial f}{\partial y} & \left( \frac{\partial f}{\partial y} \right)^2 \end{bmatrix}$$

- Block matching: same cost & model, but discretized shifts



10

11

## Regularization and MAP

- In maximum a-posteriori (MAP), the objective (or loss)  $\varepsilon$  consists of a data term (fidelity) and a prior

$$\min_{\mathbf{d}} \varepsilon_{\text{data}}(f(\mathbf{d}), g) + \varepsilon_{\text{prior}}(\mathbf{d})$$

$$\Leftrightarrow \max_{\mathbf{d}} \exp(-\varepsilon_{\text{data}}(f(\mathbf{d}), g)) \exp(-\varepsilon_{\text{prior}}(\mathbf{d}))$$

$$\Leftrightarrow \max_{\mathbf{d}} P(g|\mathbf{d})P(\mathbf{d})$$

$$\Leftrightarrow \max_{\mathbf{d}} P(\mathbf{d}|g)$$

- A common prior is a smoothness term (regularizer)



11

12

## MAP Example: KLT

- Assume a prior probability for the displacement:  $P(\mathbf{d})$  (e.g. Gaussian distribution from a motion model)
- In logarithmic domain, we now have two terms in the cost function:

$$\varepsilon(\mathbf{d}) = \sum_{\mathbf{x} \in \mathcal{N}} w(\mathbf{x}) |f(\mathbf{x} - \mathbf{d}) - g(\mathbf{x})|^2 + \lambda \|\mathbf{d} - \mathbf{d}_{\text{pred}}\|^2$$

- The standard KLT term
- A term that *drags* the solution towards the predicted displacement (cf. Kalman filtering)



12

13

## Demo: KLT

13

14

## Image Reconstruction

- Assume that  $\mathbf{u}$  is an unknown image that is observed through the linear operator  $\mathbf{G}$ :  $\mathbf{u}_0 = \mathbf{G}\mathbf{u} + \text{noise}$
- Example: blurring, linear projection
- Goal is to minimize the error  $\mathbf{u}_0 - \mathbf{G}\mathbf{u}$
- Example: squared error
- Assume that we have a prior probability for the image:  $P(\mathbf{u})$
- Example: we assume that the image should be smooth (small gradients)

14

15

## Image Reconstruction

- Minimizing

$$\varepsilon(\mathbf{u}) = \frac{1}{2}(|\mathbf{G}\mathbf{u} - \mathbf{u}_0|^2 + \lambda(|\mathbf{D}_x\mathbf{u}|^2 + |\mathbf{D}_y\mathbf{u}|^2))$$

- Gives the normal equations

$$\mathbf{G}^T\mathbf{G}\mathbf{u} - \mathbf{G}^T\mathbf{u}_0 + \lambda(\mathbf{D}_x^T\mathbf{D}_x\mathbf{u} + \mathbf{D}_y^T\mathbf{D}_y\mathbf{u}) = 0$$

- Such that

$$\mathbf{u} = (\mathbf{G}^T\mathbf{G} + \lambda(\mathbf{D}_x^T\mathbf{D}_x + \mathbf{D}_y^T\mathbf{D}_y))^{-1}\mathbf{G}^T\mathbf{u}_0$$

15

16

## Gradient Operators

- Taylor expansion of image gives

$$u(x+h, y) = u(x, y) + hu_x(x, y) + O(h^2)$$

$$u(x-h, y) = u(x, y) - hu_x(x, y) + O(h^2)$$

- Finite left/right differences give

$$\partial_x^+ u = \frac{u(x+h, y) - u(x, y)}{h} + O(h^2)$$

$$\partial_x^- u = \frac{u(x, y) - u(x-h, y)}{h} + O(h^2)$$

- Often needed: products of derivative operators

16

17

## Gradient Operators

- Squaring left (right) difference  $(\partial_x^\pm)^2 u$  gives linear error in  $h$
- Squaring central difference  $\frac{u(x+h, y) - u(x-h, y)}{2h}$  gives quadratic error in  $h$ , but leaves out every second sample

- Multiplying left and right difference

$$\partial_x^+ \partial_x^- u = \frac{u(x+h, y) - 2u(x, y) + u(x-h, y)}{h^2} = \Delta_x u$$

gives quadratic error in  $h$  (usual discrete Laplace operator)



17

18

## Demo: Image Reconstruction

- IRdemo.m



18

19

## Robust error norms

- Alternative to RANSAC (Monte-Carlo)
- Assume quadratic error: *influence* of change  $f$  to  $f+\partial f$  to the estimate is linear (why?)
- Result on set of measurements: mean
- Assume absolute error: influence of change is constant (why?)
- Result on set of measurements: median
- In general: sub-linear influence leads to robust estimates, but *non-linear*



19

20

## Smoothness

- Quadratic smoothness term: influence linear with height of edge
- Total variation smoothness (absolute value of gradient): influence constant
- With quadratic measurement error: Rudin-Osher-Fatemi (ROF) model (Physica D, 1992)

$$\min_{u \in X} \frac{\|u - g\|^2}{2\lambda} + \sum_{1 \leq i, j \leq N} |(\nabla u)_{i,j}|$$



20

21

## Total Variation (TV)

- Minimizing  $\min_{u \in X} \frac{\|u - g\|^2}{2\lambda} + \sum_{1 \leq i, j \leq N} |(\nabla u)_{i,j}|$

- Stationary point  $u - g - \lambda \operatorname{div} \left( \frac{\nabla u}{|\nabla u|} \right) = 0$

- Steepest descent

$$u^{(s+1)} = u^{(s)} - \alpha \left( u^{(s)} - g - \lambda \frac{u_{xx}u_y^2 - 2u_{xy}u_xu_y + u_{yy}u_x^2}{|\nabla u|^3} \right)$$

21

22

## Efficient TV Algorithms

- In 1D: Chambolle's algorithm (JMIV, 2004)
- In 2D:
  - Alternating direction method of multipliers (ADMM, variant of augmented Lagrangian): Split Bregman by Goldstein & Osher (SIAM 2009)
  - Based on threshold Landweber: Fast Iterative Shrinkage-Thresholding Algorithm (FISTA) by Beck & Teboulle (SIAM 2009)
  - Based on Lagrange multipliers: Primal Dual Algorithm by Chambolle & Pock (JMIV 2011)

22

23

## Demo: TV Image Denoising

23

24

## TV Image Inpainting / Convex Optimization

- Note that many problems (including quadratic and TV) are convex optimization problems
- A good first approach is to map these problems to a standard solver, e.g. CVXPY by S. Diamond and S. Boyd
- Example: minimize the total variation of an image

$$\sum_{1 \leq i, j \leq N} |(\nabla u)_{i,j}| \quad \text{under the constraint of a subset of known image values } u$$

```
prob=Problem(Minimize(tv(X)),[X[known]== MG[known]])
opt_val = prob.solve()
```

24

25

### Demo: TV Inpainting



25

27

### Non-linear LS, Dog Leg

- For comparison: LM  $r(\mathbf{x} + \delta) \approx r(\mathbf{x}) + \mathbf{J}\delta$   
 $(\mathbf{J}^T \mathbf{J} + \lambda \text{diag}(\mathbf{J}^T \mathbf{J})) \delta = \mathbf{J}^T r(\mathbf{x})$   
 $x_j \mapsto x_j + \delta_j \quad J_{ij} = \frac{\partial r_i}{\partial x_j}$

- More efficient: replace damping factor  $\lambda$  with trust region radius  $\Delta$

method	abbr.	properties
steepest descent	SD	$\delta = \mathbf{J}^T \mathbf{r}$
Gauss-Newton	GN	$\mathbf{J}^T \mathbf{J} \delta = \mathbf{J}^T \mathbf{r}$
Levenberg-Marquardt	LM	combines SD and GN by damping factor
Dog Leg	DL	combines SD and GN by trust region radius $\Delta$



27

26

### Algorithmic Taxonomy

- Minimal problems (e.g. 5 point algorithm)
  - Fully determined solution(s)
  - Analytic solvers (polynomials, Gröbner bases)
  - Numerical methods (Dogleg, Newton-Raphson)
- Overdetermined problems (e.g. OF,BA)
  - Minimization problem
  - Numerical solvers only
  - Levenberg-Marquardt (interpolation Gauss-Newton and gradient descent / trust region)

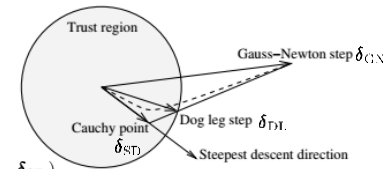


26

28

### Dog Leg (basic idea)

- initialize radius  $\Delta = 1$
- compute gain factor
- if gain factor  $> 0$



$$\mathbf{x}_{\text{new}} = \mathbf{x} + \underbrace{\delta_{\text{SD}} + \alpha(\delta_{\text{GN}} - \delta_{\text{SD}})}_{\delta_{\text{DL}}}$$

$$\|\delta_{\text{SD}}\| \leq \Delta, \quad 0 \leq \alpha \leq 1, \quad \|\delta_{\text{DL}}\| = \Delta$$

- grow/shrink  $\Delta$  and update gain factor
- if update and residual nonzero goto 3



28

29

## Optical Flow

- Minimizing (lecture 4)  $\varepsilon(\mathbf{v}_h) = \sum_{\mathcal{R}} w |\nabla^T f f_t \mathbf{v}_h|^2$
- Under the constraint  $|\mathbf{v}_h|^2 = 1$
- Using Lagrangian multiplier leads to the minimization problem
 
$$\varepsilon_T(\mathbf{v}_h, \lambda) = \varepsilon(\mathbf{v}_h) + \lambda(1 - |\mathbf{v}_h|^2)$$
- This is the *total least squares* formulation to determine the flow

29

30

## Optical Flow

- Solution is given by the eigenvalue problem

$$\left( \sum_{\mathcal{R}} w \begin{bmatrix} \nabla f \\ f_t \end{bmatrix} [\nabla^T f f_t] \right) \mathbf{v}_h = \lambda \mathbf{v}_h$$

$$\mathbf{T} \mathbf{v}_h = \lambda \mathbf{v}_h$$

- The matrix term  $\mathbf{T}$  is the spatio-temporal structure tensor
- The eigenvector with the smallest eigenvalue is the solution (up to normalization of homogeneous element)

30

## Optical Flow

- *Local* flow estimation
  - Design question:  $w$  and  $\mathcal{R}$
  - Aperture problem: motion at linear structures can only be estimated in normal direction (underdetermined)
  - Inferring limited
- *Global* flow instead



31

32

## Optical Flow

- Minimizing BCCE over the whole image with additional smoothness term

$$\varepsilon(\mathbf{f}) = \frac{1}{2} \int_{\Omega} (\langle \mathbf{f} | \nabla g \rangle + g_t)^2 + \lambda (|\nabla f_1|^2 + |\nabla f_2|^2) dx dy$$

- Gives the iterative Horn & Schunck method (details will follow in the lecture on variational methods)

$$\mathbf{f}^{(s+1)} = \bar{\mathbf{f}}^{(s)} - \frac{1}{\lambda^2 + |\nabla g|^2} (\langle \bar{\mathbf{f}}^{(s)} | \nabla g \rangle + g_t) \nabla g$$

32



33

### Graph Algorithms

- All examples so far: vectors as solutions, i.e. finite set of (pseudo) continuous values
- Now: discrete (and binary) values
- Directly related to (labeled) graph-based optimization
- In probabilistic modeling (on regular grid): Markov random fields



33

34

### Graphs

- Graph: algebraic structure  $G=(V, E)$
- Nodes  $V=\{v_1, v_2, \dots, v_n\}$
- Arcs  $E=\{e_1, e_2, \dots, e_m\}$ , where  $e_k$  is incident to
  - an unordered pair of nodes  $\{v_i, v_j\}$
  - an ordered pair of nodes  $(v_i, v_j)$  (directed graph)
  - degree of node: number of incident arcs
- Weighted graph: costs assigned to nodes or arcs

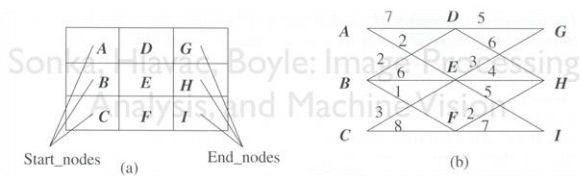


34

35

### 1D: Dynamic Programming

- Problem: find optimal path from source node  $s$  to sink node  $t$
- Principle of Optimality: If the optimal path  $s-t$  goes through  $r$ , then both  $s-r$  and  $r-t$ , are also optimal

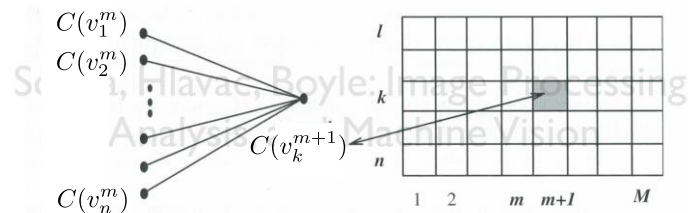


35

36

### 1D: Dynamic Programming

- $C(v_k^{m+1})$  is the new cost assigned to node  $v_k$
- $g^m(i, k)$  is the partial path cost between nodes  $v_i$  and  $v_k$



36

37

### 1D: Dynamic Programming

- $C(v_k^{m+1})$  is the new cost assigned to node  $v_k$
- $g^m(i, k)$  is the partial path cost between nodes  $v_i$  and  $v_k$

$$C(v_k^{m+1}) = \min_i (C(v_i^m) + g^m(i, k))$$

$$\min (C(v^1, v^2, \dots, v^M)) = \min_{k=1, \dots, n} (C(v_k^M))$$

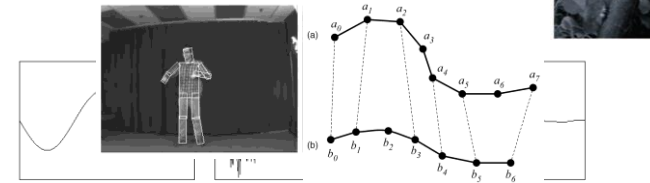


37

38

### Examples

- Shortest path computation (contours / intelligent scissors)
- 1D signal restoration (denoising)
- Tree labeling (pictorial structures)
- Matching of sequences (curves)



38

39

### Markov property

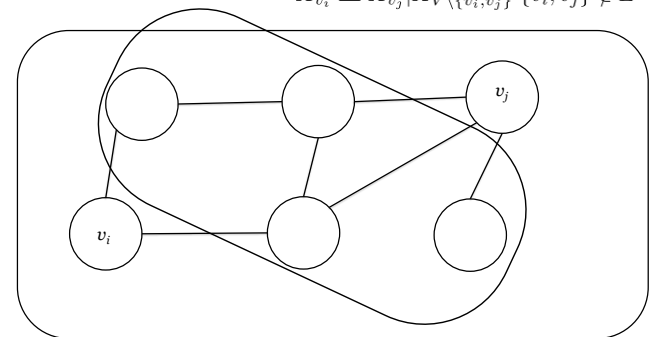
- Markov chain: memoryless process with r.v.  $X$
- Markov random field (undirected graphical model): random variables (e.g. labels) over nodes with Markov property (conditional independence)
  - Pairwise  $X_{v_i} \perp\!\!\!\perp X_{v_j} | X_{V \setminus \{v_i, v_j\}}$   $\{v_i, v_j\} \notin E$
  - Local  $X_v \perp\!\!\!\perp X_{V \setminus (\{v\} \cup N(v))} | X_{N(v)}$
  - Global  $X_A \perp\!\!\!\perp X_B | X_S$  where every path from a node in  $A$  to node in  $B$  passes through  $S$



39

### Conditional Independence

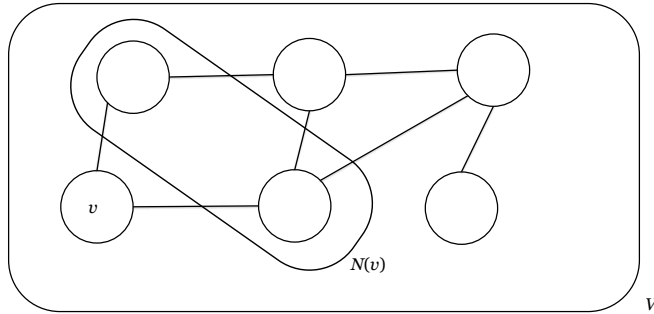
$$X_{v_i} \perp\!\!\!\perp X_{v_j} | X_{V \setminus \{v_i, v_j\}} \{v_i, v_j\} \notin E$$



40

## Conditional Independence

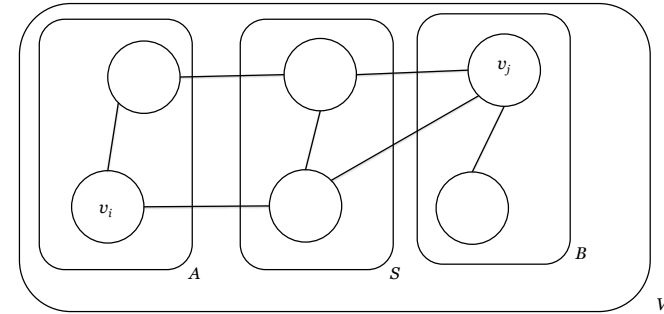
$$X_v \perp\!\!\!\perp X_{V \setminus (\{v\} \cup N(v))} \mid X_{N(v)}$$



41

## Conditional Independence

$$X_A \perp\!\!\!\perp X_B \mid X_S$$



42

## Terminology

- If joint density strictly positive: Gibbs RF
- Ising model (interacting magnetic spins), energy given as Hamiltonian function

$$\varepsilon(X_V) = - \sum_{e_k = \{v_i, v_j\} \in E} J_{e_k} X_{v_i} X_{v_j} - \sum_{v_j} h_{v_j} X_{v_j}$$

- General form

$$\varepsilon(X_V) = \lambda \sum_{e_k = \{v_i, v_j\} \in E} V(X_{v_i}, X_{v_j}) + \sum_{v_j} D(X_{v_j})$$

- Configuration probability  $P(X_V) \propto \exp(-\varepsilon(X_V))$

43

## Gibbs Model / Markov Random Field

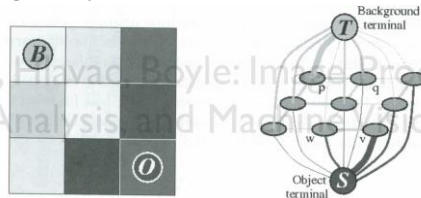
- Attempts to generalize dynamic programming to higher dimensions unsuccessful
- Minimize  $C(f) = C_{\text{data}}(f) + C_{\text{smooth}}(f)$  using arc-weighted graphs  $G_{st} = (V \cup \{s, t\}, E)$
- Two special terminal nodes, source  $s$  (e.g. object) and sink  $t$  (e.g. background) hard-linked with seed points

44

45

## Graph Cut: Two types of arcs

- n-links: connecting neighboring pixels, cost given by the smoothness term  $V$
- t-links: connecting pixels and terminals, cost given by the data term  $D$



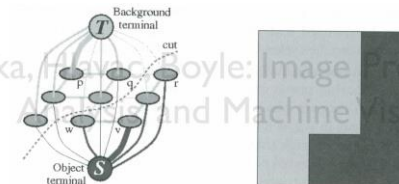
LiU LINKÖPINGS  
UNIVERSITET

45

46

## Graph Cut

- $s-t$  cut is a set of arcs, such that the nodes and the remaining arcs form two disjoint graphs with points sets  $S$  and  $T$
- cost of cut: sum of arc cost
- minimum  $s-t$  cut problem (dual: maximum flow problem)



LiU LINKÖPINGS  
UNIVERSITET

46

47

## Graph Cut

- n-link costs: large if two nodes belong to same segment, e.g. inverse gradient magnitude, Gauss, Potts model
- t-link costs:
  - $K$  for hard-linked seed points ( $K >$  maximum sum of data terms)
  - $o$  for the opposite seed point
- Submodularity  $V(\alpha, \alpha) + V(\beta, \beta) \leq V(\alpha, \beta) + V(\beta, \alpha)$

LiU LINKÖPINGS  
UNIVERSITET

47

48

## Demonstration

LiU LINKÖPINGS  
UNIVERSITET

48

49

## Examples / Discussion

- Binary problems solvable in polynomial time (albeit slow)
  - Binary image restoration
  - Bipartite matching (perfect assignment of graphs)
- N-ary problems (more than two terminals) are NP-hard and can only be approximated (e.g.  $\alpha$ -expansion move)

49



Michael Felsberg  
michael.felsberg@liu.se

[www.liu.se](http://www.liu.se)



50