

Learning with Structured Inputs and Outputs

Christoph Lampert

Institute of Science and Technology Austria
Vienna (Klosterneuburg), Austria

INRIA Summer School 2010



Institute of Science and Technology

Overview...

- 15:30–16:30 Intro + Conditional Random Fields
- 16:30–16:40 — mini-break —
- 16:40–17:30 Structured Support Vector Machines

Slides and References: <http://www.christoph-lampert.org>

"Normal" Machine Learning:

$$f : \mathcal{X} \rightarrow \mathbb{R}.$$

- inputs \mathcal{X} can be any kind of objects
- output y is a real number
 - ▶ classification, regression, density estimation, ...

Structured Output Learning:

$$f : \mathcal{X} \rightarrow \mathcal{Y}.$$

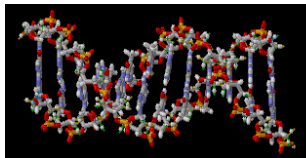
- inputs \mathcal{X} can be any kind of objects
- outputs $y \in \mathcal{Y}$ are complex (structured) objects
 - ▶ images, text, audio, folds of a protein

What is structured data?

Ad hoc definition: data that consists of several parts, and not only the parts themselves contain information, but also the way in which the parts belong together.

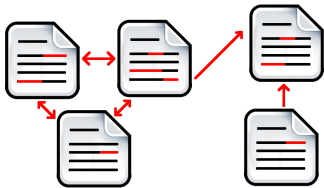
Jemand musste Josef K. verleumdet haben, denn ohne dass er etwas Böses getan hätte, wurde er eines Morgens verhaftet. »Wie ein Hund! « sagte er, es war, als sollte die Scham ihn überleben. Als Gregor Samsa eines Morgens aus unruhigen Träumen erwachte, fand er sich in seinem Bett zu einem ungeheueren Ungeziefer verwandelt. Und es war ihnen wie eine Bestätigung ihrer neuen Träume und guten Absichten, als am Ziele ihrer Fahrt die Tochter als erste sich erhob und ihren jungen Körper dehnte. »Es ist ein eigentümlicher Apparat«, sagte der Offizier zu dem Forschungsreisenden und überblickte mit einem gewissermaßen

Text

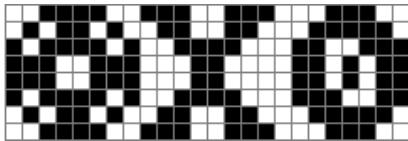


Source: wikipedia.org

Molecules / Chemical Structures



Documents/HyperText



Images

What is structured output prediction?

Ad hoc definition: predicting *structured* outputs from input data
(in contrast to predicting just a single number, like in classification or regression)

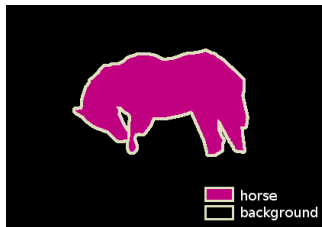
- Natural Language Processing:
 - ▶ Automatic Translation (output: sentences)
 - ▶ Sentence Parsing (output: parse trees)
- Bioinformatics:
 - ▶ Secondary Structure Prediction (output: bipartite graphs)
 - ▶ Enzyme Function Prediction (output: path in a tree)
- Speech Processing:
 - ▶ Automatic Transcription (output: sentences)
 - ▶ Text-to-Speech (output: audio signal)
- Robotics:
 - ▶ Planning (output: sequence of actions)

This lecture: Structured Prediction in Computer Vision

Computer Vision Example: Semantic Image Segmentation



input: images



output: segmentation masks

- input space $\mathcal{X} = \{images\} \hat{=} [0, 255]^{3 \cdot M \cdot N}$
- output space $\mathcal{Y} = \{segmentation\ masks\} \hat{=} \{0, 1\}^{M \cdot N}$
- (structured output) prediction function: $f : \mathcal{X} \rightarrow \mathcal{Y}$

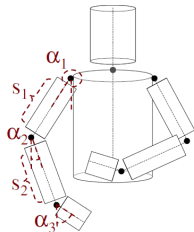
$$f(x) := \operatorname{argmin}_{y \in \mathcal{Y}} E(x, y)$$

- energy function $E(x, y) = \sum_i w_i^\top \varphi_u(x_i, y_i) + \sum_{i,j} w_{ij}^\top \varphi_p(y_i, y_j)$

Computer Vision Example: Human Pose Estimation



input: image



body model



output: model fit

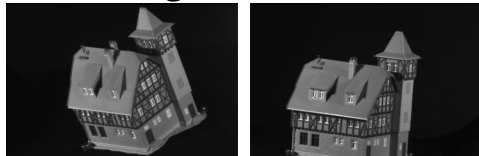
- input space $\mathcal{X} = \{images\}$
- output space $\mathcal{Y} = \{positions/angle\ of\ K\ body\ parts\} \triangleq \mathbb{R}^{3K}$.
- prediction function: $f : \mathcal{X} \rightarrow \mathcal{Y}$

$$f(x) := \operatorname{argmin}_{y \in \mathcal{Y}} E(x, y)$$

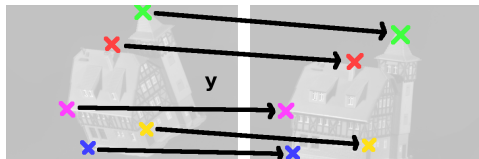
- energy $E(x, y) = \sum_i w_i^\top \varphi_{fit}(x_i, y_i) + \sum_{i,j} w_{ij}^\top \varphi_{pose}(y_i, y_j)$

Computer Vision Example: Point Matching

input: image pairs



output: mapping $y : x_i \leftrightarrow y(x_i)$



- prediction function: $f : \mathcal{X} \rightarrow \mathcal{Y}$

$$f(x) := \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y)$$

- scoring function $F(x, y) = \sum_i w_i^\top \varphi_{sim}(x_i, y(x_i)) + \sum_{i,j} w_{ij}^\top \varphi_{geom}(x_i, x_j, y(x_i), y(x_j))$

Computer Vision Example: Object Localization

input:
image



output:
object position
(left, top
right, bottom)



- input space $\mathcal{X} = \{images\}$
- output space $\mathcal{Y} = \mathbb{R}^4$ bounding box coordinates
- scoring function $F(x, y) = w^\top \varphi(x, y)$ where $\varphi(x, y) = h(x|_y)$ is a feature vector for an image region, e.g. bag-of-visual-words.
- prediction function: $f : \mathcal{X} \rightarrow \mathcal{Y}$

$$f(x) := \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y)$$

Computer Vision Examples: Summary

Image Segmentation

$$y = \operatorname{argmin}_{y \in \{0,1\}^N} E(x, y) \quad E(x, y) = \sum_i w_i^\top \varphi(x_i, y_i) + \sum_{i,j} w_{ij}^\top \varphi(y_i, y_j)$$

Pose Estimation

$$y = \operatorname{argmin}_{y \in \mathbb{R}^{3K}} E(x, y) \quad E(x, y) = \sum_i w_i^\top \varphi(x_i, y_i) + \sum_{i,j} w_{ij}^\top \varphi(y_i, y_j)$$

Point Matching

$$y = \operatorname{argmax}_{y \in \Pi_n} F(x, y) \quad F(x, y) = \sum_i w_i^\top \varphi(x_i, y_i) + \sum_{i,j} w_{ij}^\top \varphi(y_i, y_j)$$

Object Localization

$$y = \operatorname{argmax}_{y \in \mathbb{R}^4} F(x, y) \quad F(x, y) = w^\top \varphi(x, y)$$

Unified Setup

Predict structured output by maximization

$$y = \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y)$$

of a compatibility function

$$F(x, y) = \langle w, \varphi(x, y) \rangle$$

that is linear in a parameter vector w .

A generic structured prediction problem

- \mathcal{X} : arbitrary input domain
- \mathcal{Y} : structured output domain, decompose $y = (y_1, \dots, y_k)$
- Prediction function $f : \mathcal{X} \rightarrow \mathcal{Y}$ by

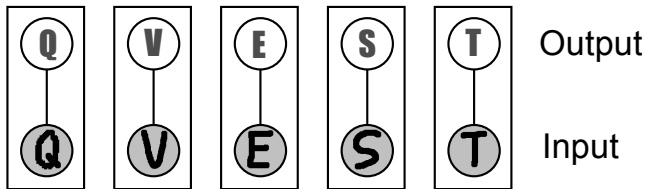
$$f(x) = \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y)$$

- Compatibility function (or negative of "energy")

$$\begin{aligned} F(x, y) &= \langle w, \varphi(x, y) \rangle \\ &= \sum_{i=1}^k w_i^\top \varphi_i(y_i, x) && \text{unary terms} \\ &\quad + \sum_{i,j=1}^k w_{ij}^\top \varphi_{ij}(y_i, y_j, x) && \text{binary terms} \\ &\quad + \dots && \text{higher order terms (sometimes)} \end{aligned}$$

Example: Sequence Prediction – Handwriting Recognition

- \mathcal{X} = 5-letter word images , $x = (x_1, \dots, x_5)$, $x_j \in \{0, 1\}^{300 \times 80}$
- \mathcal{Y} = ASCII translation , $y = (y_1, \dots, y_5)$, $y_j \in \{A, \dots, Z\}$.
- feature functions has only unary terms
 $\varphi(x, y) = (\varphi_1(x, y_1), \dots, \varphi_5(x, y_5))$.
- $F(x, y) = \langle w_1, \varphi_1(x, y_1) \rangle + \dots + \langle w_5, \varphi_5(x, y_5) \rangle$



Advantage: computing $y^* = \operatorname{argmax}_y F(x, y)$ is easy.

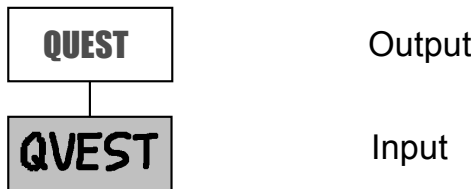
We can find each y_i^* independently, check $5 \cdot 26 = 130$ values.

Problem: only local information, we can't correct errors.

Example: Sequence Prediction – Handwriting Recognition

- \mathcal{X} = 5-letter word images , $x = (x_1, \dots, x_5)$, $x_j \in \{0, 1\}^{300 \times 80}$
- \mathcal{Y} = ASCII translation , $y = (y_1, \dots, y_5)$, $y_j \in \{A, \dots, Z\}$.
- one *global* feature function

$$\varphi(x, y) = \begin{cases} (0, \dots, 0, \overbrace{\Phi(x)}^{y\text{th pos.}}, 0, \dots, 0) & \text{if } y \in \mathcal{D} \text{ dictionary,} \\ (0, \dots, 0, 0, 0, \dots, 0) & \text{otherwise.} \end{cases}$$



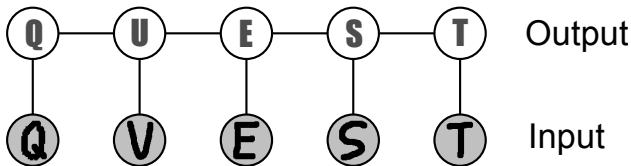
Advantage: access to global information, e.g. from dictionary \mathcal{D} .

Problem: $\operatorname{argmax}_y \langle w, \varphi(x, y) \rangle$ has to check $26^5 = 11881376$ values.
We need separate training data for each word.

Example: Sequence Prediction – Handwriting Recognition

- \mathcal{X} = 5-letter word images , $x = (x_1, \dots, x_5)$, $x_j \in \{0, 1\}^{300 \times 80}$
- \mathcal{Y} = ASCII translation , $y = (y_1, \dots, y_5)$, $y_j \in \{A, \dots, Z\}$.
- feature function with unary and pairwise terms

$$\varphi(x, y) = \left(\varphi_1(y_1, x), \varphi_2(y_2, x), \dots, \varphi_5(y_5, x), \right. \\ \left. \varphi_{1,2}(y_1, y_2), \dots, \varphi_{4,5}(y_4, y_5) \right)$$

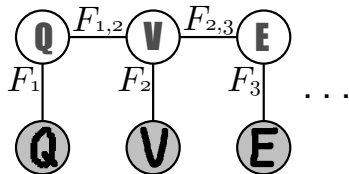


Compromise: computing y^* is still efficient (Viterbi best path)

Compromise: neighbor information allows correction of local errors.

Example: Sequence Prediction – Handwriting Recognition

- $\varphi(x, y) = (\varphi_1(y_1, x), \dots, \varphi_5(y_5, x), \varphi_{1,2}(y_1, y_2), \dots, \varphi_{4,5}(y_4, y_5))$.
- $w = (w_1, \dots, w_5, w_{1,2}, \dots, w_{4,5})$
- $F(x, y) = \langle w, \varphi(x, y) \rangle = F_1(x_1, y_1) + \dots + F_5(x_5, y_5) + F_{1,2}(y_1, y_2) + \dots, F_{4,5}(y_4, y_5)$.

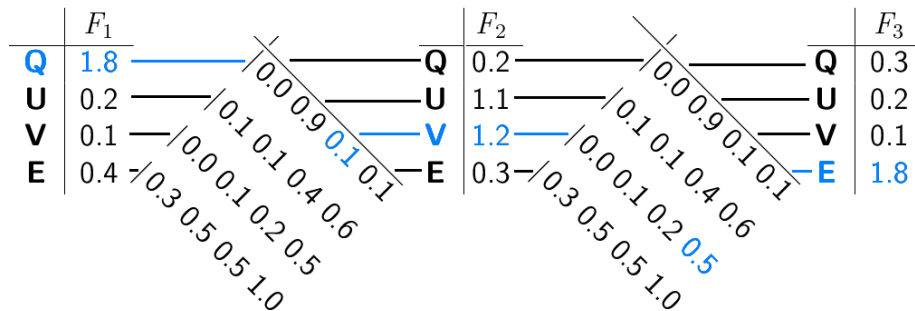


$F_{i,i+1}$	Q	U	V	E		F_1		F_2		F_3
Q	0.0	0.9	0.1	0.1	Q	1.8	Q	0.2	Q	0.3
U	0.1	0.1	0.4	0.6	U	0.2	U	1.1	U	0.2
V	0.0	0.1	0.2	0.5	V	0.1	V	1.2	V	0.1
E	0.3	0.5	0.5	1.0	E	0.4	E	0.3	E	1.8

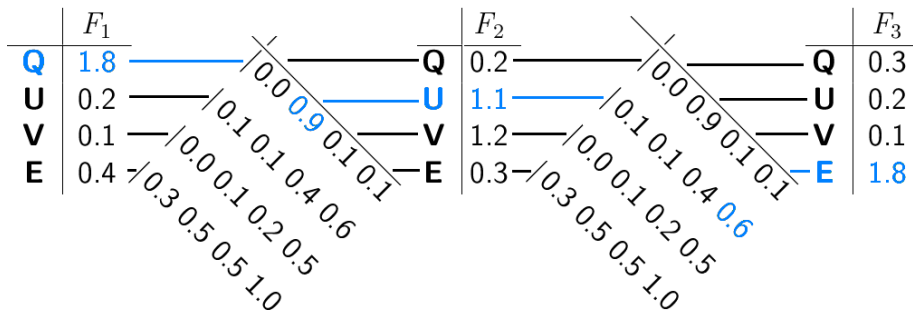
Every $y \in \mathcal{Y}$ corresponds to a path. $\operatorname{argmax}_{y \in \mathcal{Y}}$ is *best path*.

	F_1			F_2			F_3		
Q	1.8	—	/0.0	—	Q	0.2	—	Q	0.3
U	0.2	—	/0.1	—	U	1.1	—	U	0.2
V	0.1	—	/0.0	—	V	1.2	—	V	0.1
E	0.4	—	/0.3	—	E	0.3	—	E	1.8
		0.3	0.5	0.5		0.3	0.5	0.5	
		0.1	0.2	0.5		0.1	0.2	0.5	
		0.1	0.4	0.6		0.1	0.4	0.6	
		0.1	0.9	0.1		0.1	0.9	0.1	
		0.1	1.0	0.1		0.1	1.0	0.1	

$$\begin{aligned}
 F(\text{VUQ}, x) &= F_1(\text{V}, x_1) + F_{12}(\text{V}, \text{U}) + F_2(\text{U}, x_2) + F_{23}(\text{U}, \text{Q}) + F_3(\text{Q}, x_3) \\
 &= 0.1 + 0.1 + 1.1 + 0.1 + 0.3 = \mathbf{1.7}
 \end{aligned}$$



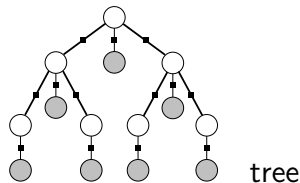
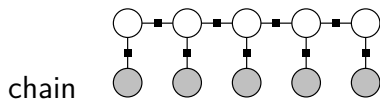
Maximal per-letter scores. Total: $1.8 + 0.1 + 1.2 + 0.5 + 1.8 = 5.4$



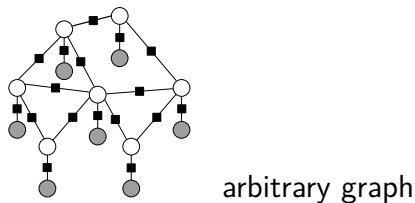
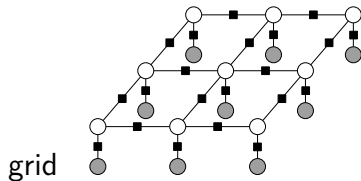
Best (=Viterbi) path. Total: $1.8 + 0.9 + 1.1 + 0.6 + 1.8 = \mathbf{6.2}$

Many popular models have **unary** and **pairwise** terms.

Yesterday's lecture: how to evaluate $\operatorname{argmax}_y F(x, y)$.



- loop-free graphs: Viterbi decoding / dynamic programming



- loopy graphs: approximate inference (e.g. loopy BP)

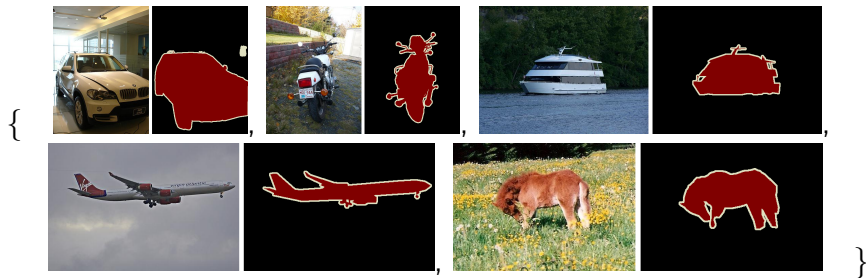
Today: how to learn a good function $F(x, y)$ from training data.

Parameter Learning in Structured Models

- Given: parametric model (family): $F(x, y) = \langle w, \varphi(x, y) \rangle$
- Given: prediction method: $f(x) = \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y)$
- Not given: parameter vector w (high-dimensional)

Supervised Training:

- Given: example pairs $\{(x^1, y^1), \dots, (x^n, y^n)\} \subset \mathcal{X} \times \mathcal{Y}$.
typical inputs with "the right" outputs for them.
- Task: determine "good" w



- Task: determine "good" w

Probabilistic Training of Structured Models

(Conditional Random Fields)

Probabilistic Models for Structured Prediction

Establish a one-to-one correspondence between

- compatibility function $F(x, y)$ and
- conditional probability distribution $p(y|x)$,

such that

$$\operatorname{argmax}_{y \in \mathcal{Y}} F(x, y) = \operatorname{argmax}_{y \in \mathcal{Y}} p(y|x)$$

= maximize the posterior probability = Bayes optimal decision

Gibbs Distribution

$$p(y|x) = \frac{1}{Z(x)} e^{F(x,y)} \quad \text{with } Z(x) = \sum_{y \in \mathcal{Y}} e^{F(x,y)}$$

$$\operatorname{argmax}_{y \in \mathcal{Y}} p(y|x) = \operatorname{argmax}_{y \in \mathcal{Y}} \frac{e^{F(x,y)}}{Z(x)} = \operatorname{argmax}_{y \in \mathcal{Y}} e^{F(x,y)} = \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y)$$



For learning, we make the dependence on w explicit:

$$F(x, y, w) = \langle w, \varphi(x, y) \rangle$$

Parametrized Gibbs Distribution – Loglinear Model

$$p(y|x, w) = \frac{1}{Z(x, w)} e^{F(x, y, w)} \quad \text{with } Z(x, w) = \sum_{y \in \mathcal{Y}} e^{F(x, y, w)}$$

Probabilistic Models for Structured Prediction

Supervised training:

- Given: i.i.d. example pairs $\mathcal{D} = \{(x^1, y^1), \dots, (x^n, y^n)\}$.
- Task: determine w .
- Idea: Maximize the posterior probability $p(w|\mathcal{D})$

Use Bayes' rule to infer posterior distribution for w :

$$\begin{aligned} p(w|\mathcal{D}) &= \frac{p(\mathcal{D}|w)p(w)}{p(\mathcal{D})} = \frac{p((x^1, y^1), \dots, (x^n, y^n)|w)p(w)}{p(\mathcal{D})} \\ &\stackrel{i.i.d.}{=} p(w) \prod_{i=1}^n \frac{p(x^i, y^i|w)}{p(x^i, y^i)} = p(w) \prod_{i=1}^n \frac{p(y^i|x^i, w)p(x^i|w)}{p(y^i|x^i)p(x^i)} \end{aligned}$$

assuming x does not depend on w : $p(x|w) \equiv p(x)$

$$= p(w) \prod_{i=1}^n \frac{p(y^i|x^i, w)}{p(y^i|x^i)}$$

Posterior probability:

$$p(w|\mathcal{D}) = p(w) \prod_i \frac{p(y^i|x^i, w)}{p(y^i|x^i)}$$

Prior Data Term (loglinear)

Typical Prior choice:

- **Gaussian** $p(w) := \text{const.} \cdot e^{-\frac{1}{2\sigma^2}\|w\|^2}$
 - ▶ Gaussianity is often a good guess if you know nothing else
 - ▶ easy to handle, analytically und numerically

Less common:

- **Laplacian** $p(w) := \text{const.} \cdot e^{-\frac{1}{\sigma}\|w\|}$
 - ▶ if we expect *sparsity* in w
 - ▶ but: properties of the optimization problem change

Learning w from $\mathcal{D} \equiv$ Maximizing posterior probability for w

$$w^* = \operatorname{argmax}_{w \in \mathbb{R}^d} p(w|\mathcal{D}) = \operatorname{argmin}_{w \in \mathbb{R}^d} \left[-\log p(w|\mathcal{D}) \right]$$

$$\begin{aligned} -\log p(w|\mathcal{D}) &= -\log \left[p(w) \prod_i \frac{p(y^i|x^i, w)}{p(y^i|x^i)} \right] \\ &= -\log p(w) - \underbrace{\sum_{i=1}^n \log p(y^i|x^i, w)}_{=\frac{1}{Z} e^{F(x^i, y^i, w)}} + \underbrace{\sum_{i=1}^n \log p(y^i|x^i)}_{\text{indep. of } w} \\ &= \frac{\|w\|^2}{2\sigma^2} - \sum_{i=1}^n \left[F(x^i, y^i, w) - \log Z(x^i, w) \right] + \text{const.} \\ &= \frac{\|w\|^2}{2\sigma^2} - \underbrace{\sum_{i=1}^n \left[\langle w, \varphi(x^i, y^i) \rangle - \log \sum_{y \in \mathcal{Y}} e^{\langle w, \varphi(x^i, y) \rangle} \right]}_{=:\mathcal{L}(w)} + \text{c.} \end{aligned}$$

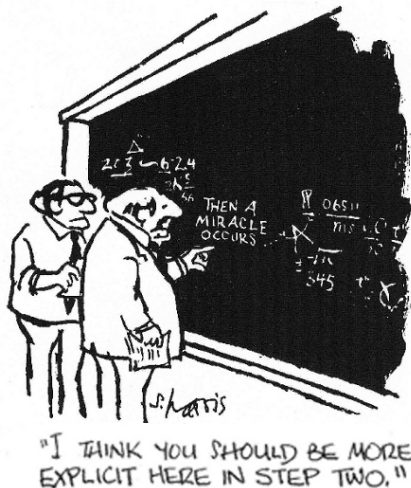
Why is everything an optimization problem?

Why all the formulas?

Why not simply teach algorithms?

Because...

- we want to separate between:
 - ▶ what is our ideal goal?
= **objective function**
 - ▶ (how) do we achieve it?
= **optimization method**
- defining a goal helps in **understanding** the problem
- mathematical formulation allows **re-using existing algorithms** (developed for different tasks)



(Negative) (Regularized) Conditional Log-Likelihood (of \mathcal{D})

$$\mathcal{L}(w) = \frac{1}{2\sigma^2} \|w\|^2 - \sum_{i=1}^n \left[\langle w, \varphi(x^i, y^i) \rangle - \log \sum_{y \in \mathcal{Y}} e^{\langle w, \varphi(x^i, y) \rangle} \right]$$

Probabilistic parameter estimation or training means solving

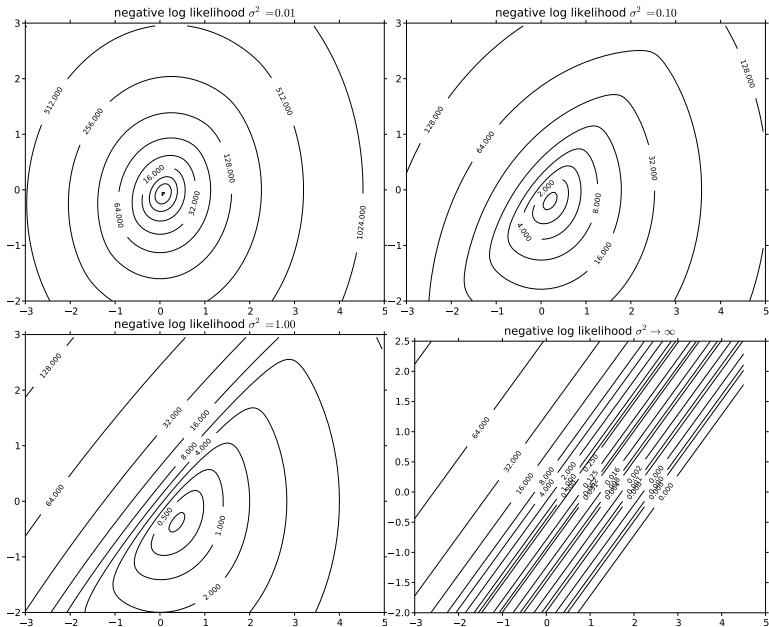
$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathcal{L}(w).$$

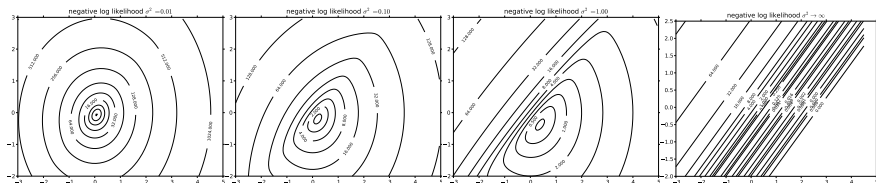
Same optimization problem as for multi-class **logistic regression**.

$$\mathcal{L}_{\text{logreg}}(w) = \frac{1}{2\sigma^2} \|w\|^2 - \sum_{i=1}^n \left[\langle w_{y^i}, \varphi(x^i) \rangle - \log \sum_{y \in \mathcal{Y}} e^{\langle w_y, \varphi(x^i) \rangle} \right]$$

with $w = (w_1, \dots, w_K)$ for K classes.

Negative Conditional Log-Likelihood for $w \in \mathbb{R}^2$:





$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathcal{L}(w).$$

- $w \in \mathbb{R}^d$ continuous (not discrete) 😊
- $w \in \mathbb{R}^d$ high-dimensional (often $d \gg 1000$) ☹️
- optimization has no side-constraints 😊
- $\mathcal{L}(w)$ differentiable 😊 → can apply gradient descent
- $\mathcal{L}(w)$ **convex!** 😊 → g.d. will find **global optimum**

Steepest Descent Minimization – minimize $\mathcal{L}(w)$

- **require:** tolerance $\epsilon > 0$
- $w_{cur} \leftarrow 0$
- **repeat**
 - ▶ $v \leftarrow -\nabla_w \mathcal{L}(w_{cur})$ descent direction
 - ▶ $\eta \leftarrow \operatorname{argmin}_{\eta \in \mathbb{R}} \mathcal{L}(w_{cur} + \eta v)$ stepsize
 - ▶ $w_{cur} \leftarrow w_{cur} + \eta v$ update
- **until** $\|v\| < \epsilon$
- **return** w_{cur}

Alternatives:

- L-BFGS (second-order descent without explicit Hessian)
- Conjugate Gradient

Summary I: Probabilistic Training (Conditional Random Fields)

- Well-defined probabilistic setting.
- $p(y|x, w)$ log-linear in $w \in \mathbb{R}^d$.
- Training: maximize $p(w|\mathcal{D})$ for data set \mathcal{D} .
- Differentiable, *convex* optimization,
- Same structure as *logistic regression*.
⇒ gradient descent will find global optimum.

For logistic regression: this is where the textbook ends. We're done.

For conditional random fields: the worst lies still ahead of us!

Problem: Computing $\nabla_w \mathcal{L}(w_{cur})$ and evaluating $\mathcal{L}(w_{cur} + \eta v)$:

$$\nabla_w \mathcal{L}(w) = \frac{1}{\sigma^2} w - \sum_{i=1}^n \left[\varphi(x^i, y^i) - \sum_{y \in \mathcal{Y}} p(y|x^i, w) \varphi(x^i, y) \right]$$

$$\mathcal{L}(w) = \frac{1}{2\sigma^2} \|w\|^2 - \sum_{i=1}^n \left[\langle w, \varphi(x^i, y^i) \rangle + \log \sum_{y \in \mathcal{Y}} e^{\langle w, \varphi(x^i, y) \rangle} \right]$$

- \mathcal{Y} typically is very (exponentially) large:
 - ▶ N parts, each with label from a set Y : $|\mathcal{Y}| = |Y|^N$,
binary image segmentation: $|\mathcal{Y}| = 2^{640 \times 480} \approx 10^{92475}$,
 - ▶ ranking N images: $|\mathcal{Y}| = N!$, e.g. $N = 1000$: $|\mathcal{Y}| \approx 10^{2568}$.

Without **structure** in \mathcal{Y} , we're lost.

$$\nabla_w \mathcal{L}(w) = \frac{1}{\sigma^2} w - \sum_{i=1}^n \left[\varphi(x^i, y^i) - \mathbb{E}_{y \sim p(y|x^i, w)} \varphi(x^i, y) \right]$$

Computing the Gradient (naive): $O(K^M nd)$

$$\mathcal{L}(w) = \frac{1}{2\sigma^2} \|w\|^2 - \sum_{i=1}^n \left[\langle w, \varphi(x^i, y^i) \rangle + \overbrace{\log \sum_{y \in \mathcal{Y}} e^{\langle w, \varphi(x^i, y) \rangle}}^{= \log Z(x^i, w)} \right]$$

Line Search (naive): $O(K^M nd)$ per evaluation of \mathcal{L}

- n : number of samples
- d : dimension of feature space
- M : number of output nodes ≈ 10 s to 1,000,000s
- K : number of possible labels of each output nodes ≈ 2 to 100s

How Does Structure in \mathcal{Y} Help?

Computing

$$Z(x, w) = \sum_{y \in \mathcal{Y}} e^{\langle w, \varphi(x, y) \rangle}$$

with φ and w decomposing over the factors

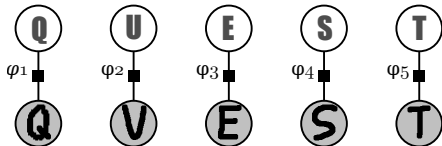
$$\varphi(x, y) = \left(\varphi_F(x_F, y_F) \right)_{F \in \mathcal{F}} \quad \text{and} \quad w = \left(w_F \right)_{F \in \mathcal{F}}$$

$$\begin{aligned} Z(x, w) &= \sum_{y \in \mathcal{Y}} e^{\sum_{F \in \mathcal{F}} \langle w_F, \varphi_F(x_F, y_F) \rangle} \\ &= \sum_{y \in \mathcal{Y}} \prod_{F \in \mathcal{F}} \underbrace{e^{\langle w_F, \varphi_F(x_F, y_F) \rangle}}_{=: \Psi_F(y_F)} \end{aligned}$$

A lot of research goes into analyzing this expression.

Case I: only unary factors

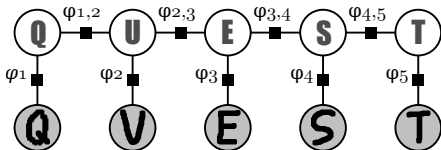
$$\mathcal{F} = \left\{ \{x_1, y_1\}, \dots, \{x_M, y_M\} \right\}.$$



$$\begin{aligned} Z &= \sum_{y \in \mathcal{Y}} \prod_{F \in \mathcal{F}} \Psi_F(y_F) \\ &= \sum_{y \in \mathcal{Y}} \prod_{i=1}^M \Psi_i(y_i) \\ &= \sum_{y_1 \in Y} \sum_{y_2 \in Y} \cdots \sum_{y_M \in Y} \Psi_1(y_1) \cdots \Psi_M(y_M) \\ &= \sum_{y_1 \in Y} \Psi_1(y_1) \sum_{y_2 \in Y} \Psi_2(y_2) \cdots \sum_{y_M \in Y} \Psi_M(y_M) \\ &= \left[\sum_{y_1 \in Y} \Psi_1(y_1) \right] \cdot \left[\sum_{y_2 \in Y} \Psi_2(y_2) \right] \cdots \left[\sum_{y_M \in Y} \Psi_M(y_M) \right] \end{aligned}$$

Case I: $O(K^M nd) \rightarrow O(MK nd)$

Case II: chain/tree with unary and pairwise factors



$$\mathcal{F} = \left\{ \{x_1, y_1\}, \{y_1, y_2\}, \dots, \{y_{M-1}, y_M\}, \{x_M, y_M\} \right\}.$$

$$\begin{aligned} Z &= \sum_{y \in \mathcal{Y}} \prod_{F \in \mathcal{F}} \Psi_F(y_F) = \sum_{y \in \mathcal{Y}} \prod_{i=1}^M \Psi_i(y_i) \prod_{i=2}^M \Psi_{i-1,i}(y_{i-1}, y_i) \\ &= \sum_{y_1 \in Y} \Psi_1 \sum_{y_2 \in Y} \Psi_{1,2} \Psi_2 \cdots \sum_{y_{M-1} \in Y} \Psi_{M-2,M-1} \Psi_{M-1} \sum_{y_M \in Y} \Psi_{M-1,M} \Psi_M \\ &= \begin{bmatrix} \square \\ \square \\ \square \\ \square \end{bmatrix}^t \cdot \begin{bmatrix} \square & \square & \square & \square \\ \square & \square & \square & \square \\ \square & \square & \square & \square \\ \square & \square & \square & \square \end{bmatrix} \cdots \begin{bmatrix} \square & \square & \square & \square \\ \square & \square & \square & \square \\ \square & \square & \square & \square \\ \square & \square & \square & \square \end{bmatrix} \cdot \begin{bmatrix} \square & \square & \square & \square \\ \square & \square & \square & \square \\ \square & \square & \square & \square \\ \square & \square & \square & \square \end{bmatrix} \begin{bmatrix} \square \\ \square \\ \square \\ \square \end{bmatrix} \end{aligned}$$

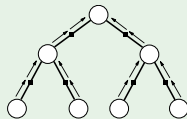
Case II: $O(\textcolor{blue}{MK}^2 nd)$

independent was $O(MKnd)$, naive $O(K^M nd)$

Message Passing in Trees

Rephrase matrix multiplication as *sending messages* from node to node:

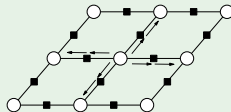
Belief Propagation (BP)



Can also compute expectations: $\mathbb{E}_{y \sim p(y|x_i, w)} \varphi(x_i, y)$.

Message Passing in Loopy Graphs

For loopy graph, send the same messages and iterate: **Loopy Belief Propagation (LBP)**



Results only in *approximation* to $Z(x, w)$, $\mathbb{E}_{y \sim p(y|x_i, w)} \varphi(x_i, y)$.
No converge guarantee, but often used in practice.

- Carsten's lecture yesterday (Wednesday)

More: ● Pawan Kumar's thesis "*Combinatorial and Convex Optimization for Probabilistic Models in Computer Vision*"
<http://ai.stanford.edu/~pawan/kumar-thesis.pdf>

$$\nabla_w \mathcal{L}(w) = \frac{1}{\sigma^2} w - \sum_{i=1}^n \left[\varphi(x^i, y^i) - \mathbb{E}_{y \sim p(y|x^i, w)} \varphi(x^i, y) \right]$$

Computing the Gradient: ~~$O(K^M nd)$~~ , $O(MK^2 nd)$ (if BP is possible):

$$\mathcal{L}(w) = \frac{1}{2\sigma^2} \|w\|^2 - \sum_{i=1}^n \left[\langle w, \varphi(x^i, y^i) \rangle + \log \sum_{y \in \mathcal{Y}} e^{\langle w, \varphi(x^i, y) \rangle} \right]$$

Line Search: ~~$O(K^M nd)$~~ , $O(MK^2 nd)$ per evaluation of \mathcal{L}

- n : number of samples $\approx 1,000$ s to $1,000,000$ s
- d : dimension of feature space
- M : number of output nodes
- K : number of possible labels of each output nodes

What, if our training set \mathcal{D} is too large (e.g. millions of examples)?

Switch to Simpler Classifier

- Train independent per-node classifiers

We lose all label dependences 😞. Still slow 😞.

Subsampling

- Create random subset $\mathcal{D}' \subset \mathcal{D}$
- Perform gradient descent to maximize $p(w|\mathcal{D}')$

Ignores all information in $\mathcal{D} \setminus \mathcal{D}'$. 😞

Parallelize

- Train several smaller models on different computers.
- Merge the models.

Follows "multi-core" trend 😊. Unclear how to merge models 😞 😊?
Doesn't reduce computation, just distributes it 😞.

What, if our training set \mathcal{D} is too large (e.g. millions of examples)?

Stochastic Gradient Descent (SGD)

- Keep maximizing $p(w|\mathcal{D})$. 😊
- In each gradient descent step:
 - ▶ Create random subset $\mathcal{D}' \subset \mathcal{D}$, ← **often just 1–3 elements!**
 - ▶ Follow approximate gradient

$$\tilde{\nabla}_w \mathcal{L}(w) = \frac{1}{\sigma^2} w - \sum_{(x^i, y^i) \in \mathcal{D}'} [\varphi(x^i, y^i) - \mathbb{E}_{y \sim p(y|x^i, w)} \varphi(x^i, y)]$$

- *Line search* no longer possible. Extra parameter: stepsize η
- SGD converges to $\operatorname{argmin}_w p(w|\mathcal{D})$! (with η chosen right)
- SGD needs more iterations, but each one is much faster

more: see L. Bottou, O. Bousquet: "*The Tradeoffs of Large Scale Learning*", NIPS 2008.
also: <http://leon.bottou.org/research/largescale>

$$\nabla_w \mathcal{L}(w) = \frac{1}{\sigma^2} w - \sum_{i=1}^n \left[\varphi(x^i, y^i) - \mathbb{E}_{y \sim p(y|x^i, w)} \varphi(x^i, y) \right]$$

Computing the Gradient: ~~$O(K^M nd)$~~ , $O(MK^2 n d)$ (if BP is possible):

$$\mathcal{L}(w) = \frac{1}{2\sigma^2} \|w\|^2 - \sum_{i=1}^n \left[\langle w, \varphi(x^i, y^i) \rangle + \log \sum_{y \in \mathcal{Y}} e^{\langle w, \varphi(x^i, y) \rangle} \right]$$

Line Search: ~~$O(K^M nd)$~~ , $O(MK^2 n d)$ per evaluation of \mathcal{L}

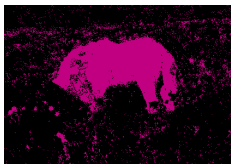
- n : number of samples
- d : dimension of feature space: $\approx \varphi_{i,j}$ 1–10s, φ_i : 100s to 10000s
- M : number of output nodes
- K : number of possible labels of each output nodes

Typical feature functions in **image segmentation**:

- $\varphi_i(y_i, x) \in \mathbb{R}^{\approx 1000}$: local image features, e.g. bag-of-words
→ $\langle w_i, \varphi_i(y_i, x) \rangle$: local classifier (like logistic-regression)
- $\varphi_{i,j}(y_i, y_j) = \mathbb{I}[y_i = y_j] \in \mathbb{R}^1$: test for same label
→ $\langle w_{ij}, \varphi_{ij}(y_i, y_j) \rangle$: penalizer for label changes (if $w_{ij} > 0$)
- combined: $\operatorname{argmax}_y p(y|x)$ is smoothed version of local cues



original



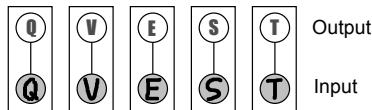
local classification



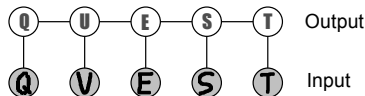
local + smoothness

Typical feature functions in **handwriting recognition**:

- $\varphi_i(y_i, x) \in \mathbb{R}^{\approx 1000}$: image representation (pixels, gradients)
→ $\langle w_i, \varphi_i(y_i, x) \rangle$: local classifier if x_i is letter y_i
- $\varphi_{i,j}(y_i, y_j) = e_{y_i} \otimes e_{y_j} \in \mathbb{R}^{26 \cdot 26}$: letter/letter indicator
→ $\langle w_{ij}, \varphi_{ij}(y_i, y_j) \rangle$: encourage/suppress letter combinations
- combined: $\operatorname{argmax}_y p(y|x)$ is "corrected" version of local cues



local classification



local + "correction"

Typical feature functions in **pose estimation**:

- $\varphi_i(y_i, x) \in \mathbb{R}^{\approx 1000}$: local image representation, e.g. HoG
→ $\langle w_i, \varphi_i(y_i, x) \rangle$: local confidence map
- $\varphi_{i,j}(y_i, y_j) = \text{good_fit}(y_i, y_j) \in \mathbb{R}^1$: test for geometric fit
→ $\langle w_{ij}, \varphi_{ij}(y_i, y_j) \rangle$: penalizer for unrealistic poses
- together: $\operatorname{argmax}_y p(y|x)$ is sanitized version of local cues



original



local classification



local + geometry

Typical feature functions for **CRFs in Computer Vision**:

- $\varphi_i(y_i, x)$: local representation, high-dimensional
→ $\langle w_i, \varphi_i(y_i, x) \rangle$: local classifier
- $\varphi_{i,j}(y_i, y_j)$: prior knowledge, low-dimensional
→ $\langle w_{ij}, \varphi_{ij}(y_i, y_j) \rangle$: penalize outliers
- learning adjusts parameters:
 - ▶ unary w_i : learn local classifiers and their importance
 - ▶ pairwise w_{ij} : learn importance of smoothing/penalization
- $\operatorname{argmax}_y p(y|x)$ is cleaned up version of local prediction

Idea: split learning of unary potentials into two parts:

- local classifiers,
- their importance.

Two-Stage Training

- pre-train $f_i^y(x) \triangleq \log p(y_i|x)$
- use $\varphi_i(y_i, x) := f_i^y(x) \in \mathbb{R}^K$ (low-dimensional)
- keep $\varphi_{ij}(y_i, y_j)$ as before
- perform CRF learning with φ_i and φ_{ij}

Advantage:

- CRF feature space much lower dimensional \rightarrow faster training
- $f_i^y(x)$ can be stronger classifiers, e.g. non-linear SVMs

Disadvantage:

- if local classifiers are really bad, CRF training cannot fix that.

Summary II — Numeric Training

(All?) Current CRF training methods rely on gradient descent.

$$\tilde{\nabla}_w \mathcal{L}(w) = \frac{w}{\sigma^2} - \sum_{i=1}^n \left[\varphi(x^i, y^i) - \sum_{y \in \mathcal{Y}} p(y|x^i, w) \varphi(x^i, y) \right] \in \mathbb{R}^d$$

The faster we can do it, the better (more realistic) models we can use.

A lot of research on speeding up CRF training:

problem	"solution"	method(s)
$ \mathcal{Y} $ too large	exploit structure <i>smart sampling</i>	(loopy) belief propagation <i>contrastive divergence</i>
n too large	mini-batches	stochastic gradient descent
d too large	discriminative φ_{unary}	two-stage training

Conditional Random Fields Summary

CRFs model *conditional probability distribution* $p(y|x) = \frac{1}{Z} e^{\langle w, \varphi(x,y) \rangle}$.

- $x \in \mathcal{X}$ is **arbitrary input**
 - ▶ we need a feature function $\varphi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$,
 - ▶ we don't need to know/estimate $p(x)$!
- $y \in \mathcal{Y}$ is **(almost arbitrary) output**
 - ▶ we need structure in \mathcal{Y} to make training/inference feasible
 - ▶ e.g. graph labeling: local properties plus neighborhood relations
- $w \in \mathbb{R}^d$ is **parameter vector**
 - ▶ *CRF training* is solving $w = \operatorname{argmax}_w p(w|\mathcal{D})$ for training set \mathcal{D} .
- **after learning, we can predict:** $f(x) := \operatorname{argmax}_y p(y|x)$.

Many many many more details in Machine Learning literature, e.g.:

- [M. Wainwright, M. Jordan: "Graphical Models, Exponential Families, and Variational Inference", now publishers, 2008]. Free PDF at <http://www.nowpublishers.com>
- [D. Koller, N. Friedman: "Probabilistic Graphical Models", MIT Press, 2010]
60 EUR (i.e. less than 5 cents per page)

Maximum Margin Training of Structured Models (Structure Output SVMs)

Example: Object Localization

Learn a **localization function**:

$$f : \mathcal{X} \rightarrow \mathcal{Y} \quad \text{with } \mathcal{X} = \{\text{images}\}, \mathcal{Y} = \{\text{bounding boxes}\}.$$



Compatibility function:

$$F : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R} \quad \text{with} \quad F(x, y) = \langle w, \varphi(x, y) \rangle$$

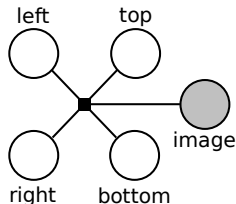
with $\varphi(x, y)$, e.g., the BoW histogram of the region y in x .

$$f(x) = \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y)$$

Learn a **localization** function:

$$f(x) = \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y) \quad \text{with } F(x, y) = \langle w, \varphi(x, y) \rangle.$$

$\varphi(x, y) = \text{BoW histogram of region } y = [\text{left}, \text{top}, \text{right}, \text{bottom}] \text{ in } x.$



$$p(y|x, w) = \frac{1}{Z} \exp \left(\langle w, \varphi(x, y) \rangle \right) \quad \text{with } Z = \sum_{y' \in \mathcal{Y}} \exp \left(\langle w, \varphi(x, y') \rangle \right)$$

We can't train probabilistically: φ doesn't decompose, Z too hard.
We **can do** MAP prediction: sliding window or branch-and-bound.

Can we learn structured prediction functions without " Z "?

Reminder: general form of regularized learning

Find w for decision function $f = \langle w, \varphi(x, y) \rangle$ by

$$\min_{w \in \mathbb{R}^d} \quad \lambda \|f\|^2 \quad + \quad \sum_{i=1}^n \ell(y^i, f(x^i))$$

Regularization + Error on data

Probabilistic Training: logistic loss

$$\ell(y^i, f(x^i)) := \log \sum_{y \in \mathcal{Y}} \exp \left(\langle w, \varphi(x^i, y) \rangle - \langle w, \varphi(x^i, y^i) \rangle \right)$$

Maximum-Margin Training: Hinge-loss

$$\ell(y^i, f(x^i)) := \max \left\{ 0, \max_{y \in \mathcal{Y}} \Delta(y^i, y) + \langle w, \varphi(x^i, y) \rangle - \langle w, \varphi(x^i, y^i) \rangle \right\}$$

$\Delta(y, y') \geq 0$ with $\Delta(y, y) = 0$ is *loss function* between outputs.

Conditional Random Fields:

$$\min_{w \in \mathbb{R}^d} \frac{\|w\|^2}{2\sigma^2} + \log \sum_{y \in \mathcal{Y}} \exp \left(\langle w, \varphi(x^i, y) \rangle - \langle w, \varphi(x^i, y^i) \rangle \right)$$

Unconstrained optimization, convex, differentiable objective.

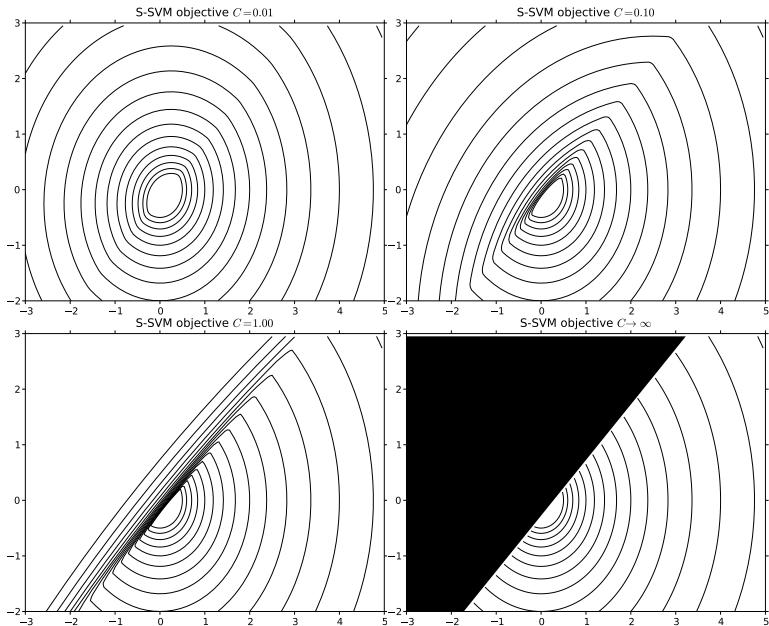
Structured Output Support Vector Machine:

$$\min_{w \in \mathbb{R}^d} \frac{1}{2} \|w\|^2 + \frac{C}{n} \sum_{i=1}^n \left[\max_{y \in \mathcal{Y}} \Delta(y^i, y) + \langle w, \varphi(x^i, y) \rangle - \langle w, \varphi(x^i, y^i) \rangle \right]_+$$

where $[t]_+ := \max\{0, t\}$.

Unconstrained optimization, convex, non-differentiable objective.

S-SVM Objective Function for $w \in \mathbb{R}^2$:



Structured Output SVM (equivalent formulation):

$$\min_{w \in \mathbb{R}^d, \xi \in \mathbb{R}_+^n} \quad \frac{1}{2} \|w\|^2 + \frac{C}{n} \sum_{i=1}^n \xi^i$$

subject to, for $i = 1, \dots, n$,

$$\max_{y \in \mathcal{Y}} \left[\Delta(y^i, y) + \langle w, \varphi(x^i, y) \rangle - \langle w, \varphi(x^i, y^i) \rangle \right] \leq \xi^i$$

n non-linear constraints, convex, differentiable objective.

Structured Output SVM (also equivalent formulation):

$$\min_{w \in \mathbb{R}^d, \xi \in \mathbb{R}_+^n} \quad \frac{1}{2} \|w\|^2 + \frac{C}{n} \sum_{i=1}^n \xi^i$$

subject to, for $i = 1, \dots, n$,

$$\Delta(y^i, y) + \langle w, \varphi(x^i, y) \rangle - \langle w, \varphi(x^i, y^i) \rangle \leq \xi^i, \quad \text{for all } y \in \mathcal{Y}$$

$|\mathcal{Y}|n$ linear constraints, convex, differentiable objective.

Example: A "True" Multiclass SVM

- $\mathcal{Y} = \{1, 2, \dots, K\}$, $\Delta(y, y') = \begin{cases} 1 & \text{for } y \neq y' \\ 0 & \text{otherwise.} \end{cases}$
- $\varphi(x, y) = \left(\mathbb{I}[y = 1]\Phi(x), \mathbb{I}[y = 2]\Phi(x), \dots, \mathbb{I}[y = K]\Phi(x) \right)$
 $= \Phi(x)e_y^\top$ with $e_y = y$ -th unit vector

Solve:

$$\min_{w, \xi} \frac{1}{2} \|w\|^2 + \frac{C}{n} \sum_{i=1}^n \xi^i$$

subject to, for $i = 1, \dots, n$,

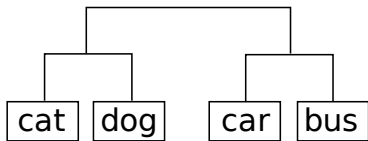
$$\langle w, \varphi(x^i, y^i) \rangle - \langle w, \varphi(x^i, y) \rangle \geq 1 - \xi^i \quad \text{for all } y \in \mathcal{Y} \setminus \{y^i\}.$$

Classification: MAP $f(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \langle w, \varphi(x, y) \rangle$

Crammer-Singer Multiclass SVM

Hierarchical Multiclass Classification

Loss function can reflect hierarchy:



$$\Delta(y, y') := \frac{1}{2}(\text{distance in tree})$$

$$\Delta(\text{cat}, \text{cat}) = 0, \quad \Delta(\text{cat}, \text{dog}) = 1, \quad \Delta(\text{cat}, \text{bus}) = 2, \quad \text{etc.}$$

Solve:

$$\min_{w, \xi} \frac{1}{2} \|w\|^2 + \frac{C}{n} \sum_{i=1}^n \xi^i$$

subject to, for $i = 1, \dots, n$,

$$\langle w, \varphi(x^i, y^i) \rangle - \langle w, \varphi(x^i, y) \rangle \geq \Delta(y^i, y) - \xi^i \quad \text{for all } y \in \mathcal{Y} \setminus \{y^i\}.$$

How to solve?

$$\min_{w \in \mathbb{R}^d} \frac{1}{2} \|w\|^2 + \frac{C}{n} \sum_{i=1}^n \left[\max_{y \in \mathcal{Y}} \Delta(y^i, y) + \langle w, \varphi(x^i, y) \rangle - \langle w, \varphi(x^i, y^i) \rangle \right]_+$$

with $[t]_+ := \max\{0, t\}$.

Unconstrained, convex, non-differentiable:

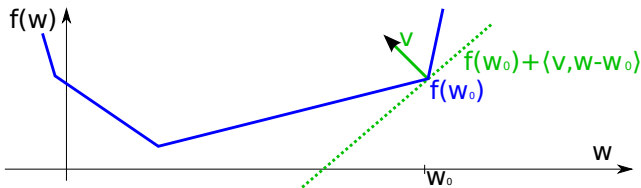
- unconstrained ☺
- convex ☺
- non-differentiable ☹ \rightarrow we can't use gradient descent directly.

Idea: use **sub-gradient descent**.

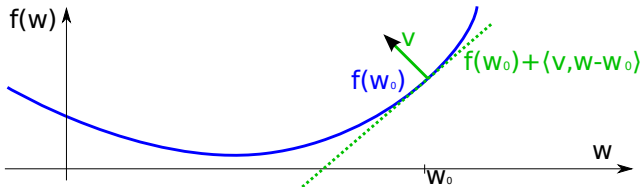
Sub-Gradients

Definition: Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function. A vector $v \in \mathbb{R}^d$ is called a **sub-gradient** of f at w_0 , if

$$f(w) \geq f(w_0) + \langle v, w - w_0 \rangle \quad \text{for all } w.$$



For differentiable f , the gradient $v = \nabla f(w_0)$ is the only subgradient.

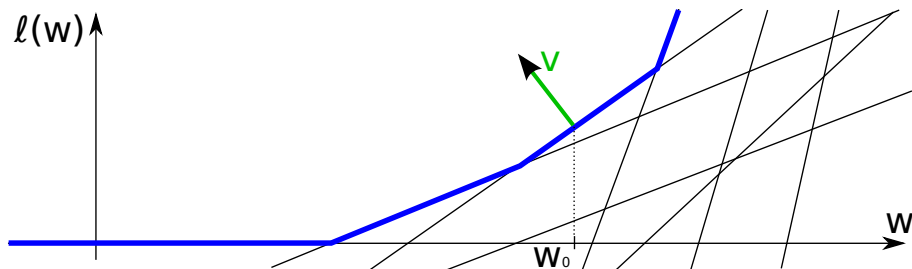


Computing a subgradient:

$$\min_w \frac{1}{2} \|w\|^2 + \frac{C}{n} \sum_{i=1}^n \ell^i(w)$$

with $\ell^i(w) = \max\{0, \max_y \ell_y^i(w)\}$, and

$$\ell_y^i(w) := \Delta(y^i, y) + \langle w, \varphi(x^i, y) \rangle - \langle w, \varphi(x^i, y^i) \rangle$$



Subgradient of ℓ^i at w_0 : find maximal (active) y , use $v = \nabla \ell_y^i(w_0)$.

Subgradient Descent S-SVM Training

input training pairs $\{(x^1, y^1), \dots, (x^n, y^n)\} \subset \mathcal{X} \times \mathcal{Y}$,

input feature map $\varphi(x, y)$, loss function $\Delta(y, y')$, regularizer C ,

input number of iterations T , stepsizes η_t for $t = 1, \dots, T$

1: $w \leftarrow \vec{0}$

2: **for** $t=1, \dots, T$ **do**

3: **for** $i=1, \dots, n$ **do**

4: $\hat{y} \leftarrow \operatorname{argmax}_{y \in \mathcal{Y}} \Delta(y^i, y) + \langle w, \varphi(x^i, y) \rangle - \langle w, \varphi(x^i, y^i) \rangle$

5: $v^i \leftarrow \varphi(x^i, \hat{y}) - \varphi(x^i, y^i)$

6: **end for**

7: $w \leftarrow w - \eta_t(w - \frac{C}{n} \sum_i v^i)$

8: **end for**

output prediction function $f(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \langle w, \varphi(x, y) \rangle$.

Obs: each update of w needs 1 MAP-like prediction per example.

We can use the same tricks as for CRFs, e.g. **stochastic updates**:

Stochastic Subgradient Descent S-SVM Training

input training pairs $\{(x^1, y^1), \dots, (x^n, y^n)\} \subset \mathcal{X} \times \mathcal{Y}$,

input feature map $\varphi(x, y)$, loss function $\Delta(y, y')$, regularizer C ,

input number of iterations T , stepsizes η_t for $t = 1, \dots, T$

1: $w \leftarrow \vec{0}$

2: **for** $t=1, \dots, T$ **do**

3: $(x^i, y^i) \leftarrow$ randomly chosen training example pair

4: $\hat{y} \leftarrow \operatorname{argmax}_{y \in \mathcal{Y}} \Delta(y^i, y) + \langle w, \varphi(x^i, y) \rangle - \langle w, \varphi(x^i, y^i) \rangle$

5: $v^i \leftarrow \varphi(x^i, \hat{y}) - \varphi(x^i, y^i)$

6: $w \leftarrow w - \eta_t(w - Cv^i)$

7: **end for**

output prediction function $f(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \langle w, \varphi(x, y) \rangle$.

Obs: each update of w needs only 1 MAP-like prediction
(but we'll need more iterations until convergence)

Numeric Solution: Solving an S-SVM like a linear SVM.

Other, equivalent, S-SVM formulation was

$$\min_{w \in \mathbb{R}^d, \xi \in \mathbb{R}_+^n} \|w\|^2 + \frac{C}{n} \sum_{i=1}^n \xi^i$$

subject to, for $i = 1, \dots, n$,

$$\langle w, \varphi(x^i, y^i) \rangle - \langle w, \varphi(x^i, y) \rangle \geq \Delta(y^i, y) - \xi^i, \quad \text{for all } y \in \mathcal{Y} \setminus \{y^i\}.$$

With feature vectors $\delta\varphi(x^i, y^i, y) := \varphi(x^i, y^i) - \varphi(x^i, y)$, this is

$$\min_{w \in \mathbb{R}^d, \xi \in \mathbb{R}_+^n} \|w\|^2 + \frac{C}{n} \sum_{i=1}^n \xi^i$$

subject to, for $i = 1, \dots, n$, for all $y \in \mathcal{Y} \setminus \{y^i\}$,

$$\langle w, \delta\varphi(x^i, y^i, y) \rangle \geq \Delta(y^i, y) - \xi^i.$$

Solve

$$\min_{w \in \mathbb{R}^d, \xi \in \mathbb{R}_+^n} \|w\|^2 + \frac{C}{n} \sum_{i=1}^n \xi^i \quad (\text{QP})$$

subject to, for $i = 1, \dots, n$, for all $y \in \mathcal{Y} \setminus \{y^i\}$,

$$\langle w, \delta\varphi(x^i, y^i, y) \rangle \geq \Delta(y^i, y) - \xi^i.$$

This has the same structure as an ordinary SVM!

- quadratic objective ☺
- linear constraints ☺

Question: Can't we use a ordinary QP solver to solve this?

Answer: Almost! We could, if there weren't $n(|\mathcal{Y}| - 1)$ constraints.

- E.g. 100 binary 16×16 images: 10^{79} constraints ☹

Solution: working set training

- It's enough if we enforce the **active constraints**.
The others will be fulfilled automatically.
- We don't know which ones are active for the optimal solution.
- But it's likely to be only a small number ← can of course be formalized.

Keep a set of potentially active constraints and update it iteratively:

- Start with working set $S = \emptyset$ (no constraints)
- Repeat until convergence:
 - ▶ Solve S-SVM training problem with constraints from S
 - ▶ Check, if solution violates any of the *full* constraint set
 - ★ if no: we found the optimal solution, *terminate*.
 - ★ if yes: add most violated constraints to S , *iterate*.

Good *practical performance* and *theoretic guarantees*:

- polynomial time convergence ϵ -close to the global optimum

Working Set S-SVM Training

input training pairs $\{(x^1, y^1), \dots, (x^n, y^n)\} \subset \mathcal{X} \times \mathcal{Y}$,

input feature map $\varphi(x, y)$, loss function $\Delta(y, y')$, regularizer C

1: $S \leftarrow \emptyset$

2: **repeat**

3: $(w, \xi) \leftarrow$ *solution to QP only with constraints from S*

4: **for** $i=1, \dots, n$ **do**

5: $\hat{y} \leftarrow \operatorname{argmax}_{y \in \mathcal{Y}} \Delta(y^i, y) + \langle w, \varphi(x^i, y) \rangle$

6: **if** $\hat{y} \neq y^i$ **then**

7: $S \leftarrow S \cup \{(x^i, \hat{y})\}$

8: **end if**

9: **end for**

10: **until** S doesn't change anymore.

output prediction function $f(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \langle w, \varphi(x, y) \rangle$.

Obs: each update of w needs 1 MAP-like prediction per example.
(but we solve globally for w , not locally)

Specialized modifications are possible: In object localization solving

$$\operatorname{argmax}_{i=1,\dots,n, y \in \mathcal{Y}} \ell_y^i(w)$$

can be done much faster than by exhaustive evaluation of

$$\operatorname{argmax}_{y \in \mathcal{Y}} \ell_y^1(w), \dots, \operatorname{argmax}_y \ell_y^n(w).$$

Single Element Working Set S-SVM Training

input ...

- 1: $S \leftarrow \emptyset$
- 2: **repeat**
- 3: $(w, \xi) \leftarrow \text{solution to QP only with constraints from } S$
- 4: $(\hat{i}, \hat{y}) \leftarrow \operatorname{argmax}_{i,y} \Delta(y^i, y) + \langle w, \varphi(x^i, y) \rangle$
- 5: **if** $\hat{y} \neq y^{\hat{i}}$ **then**
- 6: $S \leftarrow S \cup \{(x^{\hat{i}}, \hat{y})\}$
- 7: **end if**
- 8: **until** S doesn't change anymore.

Obs: each update of w needs only 1 MAP-like prediction.

"One-Slack" Structured Output SVM:

(equivalent to ordinary S-SVM by $\xi = \frac{1}{n} \sum_i \xi^i$)

$$\min_{w \in \mathbb{R}^d, \xi \in \mathbb{R}_+} \quad \frac{1}{2} \|w\|^2 + C\xi$$

subject to

$$\sum_{i=1}^n \left[\Delta(y^i, \hat{y}^i) + \langle w, \varphi(x^i, \hat{y}^i) \rangle - \langle w, \varphi(x^i, y^i) \rangle \right] \leq n\xi,$$

subject to, for all $(\hat{y}^1, \dots, \hat{y}^n) \in \mathcal{Y} \times \dots \times \mathcal{Y}$,

$$\sum_{i=1}^n \Delta(y^i, \hat{y}^i) + \langle w, \sum_{i=1}^n [\varphi(x^i, \hat{y}^i) - \varphi(x^i, y^i)] \rangle \leq n\xi,$$

$|\mathcal{Y}|^n$ linear constraints, convex, differentiable objective.

We have blown up the constraint set even further:

- 100 binary 16×16 images: 10^{177} constraints (instead of 10^{79}).

Working Set One-Slack S-SVM Training

input training pairs $\{(x^1, y^1), \dots, (x^n, y^n)\} \subset \mathcal{X} \times \mathcal{Y}$,

input feature map $\varphi(x, y)$, loss function $\Delta(y, y')$, regularizer C

1: $S \leftarrow \emptyset$

2: **repeat**

3: $(w, \xi) \leftarrow$ *solution to QP only with constraints from S*

4: **for** $i=1, \dots, n$ **do**

5: $\hat{y}^i \leftarrow \operatorname{argmax}_{y \in \mathcal{Y}} \Delta(y^i, y) + \langle w, \varphi(x^i, y) \rangle$

6: **end for**

7: $S \leftarrow S \cup \{(x^1, \dots, x^n), (\hat{y}^1, \dots, \hat{y}^n)\}$

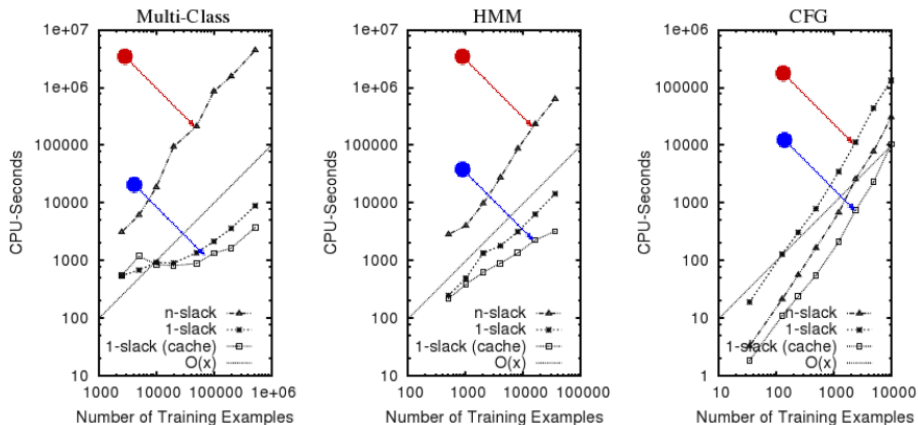
8: **until** S doesn't change anymore.

output prediction function $f(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \langle w, \varphi(x, y) \rangle$.

Often faster convergence:

We add one strong constraint per iteration instead of n weak ones.

Training Times **ordinary S-SVM** versus **one-slack S-SVM**



Observation 1: One-slack training is usually faster than n -slack.

Observation 2: Training structured models is nevertheless **slow**.

Numeric Solution: Solving an S-SVM like a non-linear SVM.

Dual S-SVM problem

$$\max_{\alpha \in \mathbb{R}_+^{n|\mathcal{Y}|}} \sum_{\substack{i=1, \dots, n \\ y \in \mathcal{Y}}} \alpha_{iy} \Delta(y^i, y) - \frac{1}{2} \sum_{\substack{y, y' \in \mathcal{Y} \\ i, i'=1, \dots, n}} \alpha_{iy} \alpha_{i'y'} \langle \delta\varphi(x^i, y^i, y), \delta\varphi(x^{i'}, y^{i'}, y') \rangle$$

subject to, for $i = 1, \dots, n$,

$$\sum_{y \in \mathcal{Y}} \alpha_{iy} \leq \frac{C}{n}.$$

Like ordinary SVMs, we can *dualize* the S-SVMs optimization:

- min becomes max,
- original (primal) variables w, ξ disappear,
- we get one dual variable α_{iy} for each primal constraint.

n linear constraints, convex, differentiable objective, $n|\mathcal{Y}|$ **variables**.

Data appears only inside inner products: **kernelize**

- Define joint kernel function $k : (\mathcal{X} \times \mathcal{Y}) \times (\mathcal{X} \times \mathcal{Y}) \rightarrow \mathbb{R}$

$$k((x, y), (x', y')) = \langle \varphi(x, y), \varphi(x', y') \rangle.$$

- k measure similarity between two *(input,output)*-pairs.
- We can express the optimization in terms of k :

$$\begin{aligned} & \langle \delta\varphi(x^i, y^i, y), \delta\varphi(x^{i'}, y^{i'}, y') \rangle \\ &= \left\langle \varphi(x^i, y^i) - \varphi(x^i, y), \varphi(x^{i'}, y^{i'}) - \varphi(x^{i'}, y') \right\rangle \\ &= \langle \varphi(x^i, y^i), \varphi(x^{i'}, y^{i'}) \rangle - \langle \varphi(x^i, y^i), \varphi(x^{i'}, y') \rangle \\ &\quad - \langle \varphi(x^i, y), \varphi(x^{i'}, y^{i'}) \rangle + \langle \varphi(x^i, y), \varphi(x^{i'}, y') \rangle \\ &= k((x^i, y^i), (x^{i'}, y^{i'})) - k((x^i, y^i), \varphi(x^{i'}, y')) \\ &\quad - k((x^i, y), (x^{i'}, y^{i'})) + k((x^i, y), \varphi(x^{i'}, y')) \\ &=: K_{ii'yy'} \end{aligned}$$

Kernelized S-SVM problem:

$$\max_{\alpha \in \mathbb{R}_+^{n|\mathcal{Y}|}} \sum_{\substack{i=1,\dots,n \\ y \in \mathcal{Y}}} \alpha_{iy} \Delta(y^i, y) - \frac{1}{2} \sum_{\substack{y, y' \in \mathcal{Y} \\ i, i'=1,\dots,n}} \alpha_{iy} \alpha_{i'y'} K_{ii'yy'}$$

subject to, for $i = 1, \dots, n$,

$$\sum_{y \in \mathcal{Y}} \alpha_{iy} \leq \frac{C}{n}.$$

- too many variables: train with **working set** of α_{iy} .

Kernelized prediction function:

$$f(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \sum_{iy'} \alpha_{iy'} k((x^i, y^i), (x, y))$$

Be careful in kernel design:

Evaluating f can easily become very slow / completely infeasible.

What do "joint kernel functions" look like?

$$k((x, y), (x', y')) = \langle \varphi(x, y), \varphi(x', y') \rangle.$$

Graphical Model: φ decomposes into components per factor:

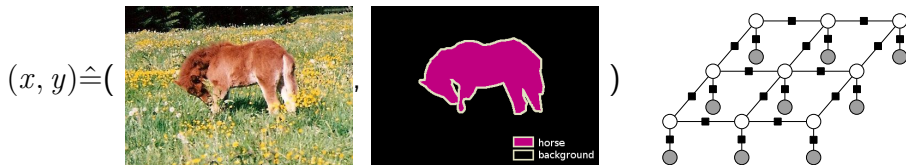
- $\varphi(x, y) = \left(\varphi_f(x_f, y_f) \right)_{f \in \mathcal{F}}$

Kernel k decomposes into sum over *factors*:

$$\begin{aligned} k((x, y), (x', y')) &= \left\langle \left(\varphi_f(x_f, y_f) \right)_{f \in \mathcal{F}}, \left(\varphi_f(x'_f, y'_f) \right)_{f \in \mathcal{F}} \right\rangle \\ &= \sum_{f \in \mathcal{F}} \langle \varphi_f(x_f, y_f), \varphi_f(x'_f, y'_f) \rangle \\ &= \sum_{f \in \mathcal{F}} k_f((x_f, y_f), (x'_f, y'_f)) \end{aligned}$$

We can define kernels for each factor (e.g. nonlinear).

Example: figure-ground segmentation with grid structure



Typical kernels: arbitrary in x , linear (or at least simple) w.r.t. y :

- Unary factors:

$$k_p((x, y_p), (x', y'_p)) = k(x_{N(p)}, x'_{N(p)}) \mathbb{I}[y_p = y'_p]$$

with k image kernel on local neighborhood $N(p)$, e.g. χ^2 or *histogram intersection*

- Pairwise factors:

$$k_{pq}((y_p, y_q), (y'_p, y'_q)) = \mathbb{I}[y_p = y'_p] \mathbb{I}[y_q = y'_q]$$

More powerful than all-linear, and MAP prediction still possible.

Example: object localization



Only one factor that includes all x and y :

$$k((x, y), (x', y')) = k_{image}(x|_y, x'|_{y'})$$

with k_{image} image kernel and $x|_y$ is image region within box y .

MAP prediction is just object localization with k_{image} -SVM.

$$f(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \sum_{iy'} \alpha_{iy'} k_{image}(x^i|_{y'}, x|_y)$$

Summary I – S-SVM Learning

Task: parameter learning

- given training set $\{(x^1, y^1), \dots, (x^n, y^n)\} \subset \mathcal{X} \times \mathcal{Y}$, learn parameters w for prediction function

$$F(x, y) = \langle w, \varphi(x, y) \rangle.$$

- predict $f : \mathcal{X} \rightarrow \mathcal{Y}$ by $f(x) = \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y)$ ($\hat{=}$ MAP)

Solution from maximum-margin principle:

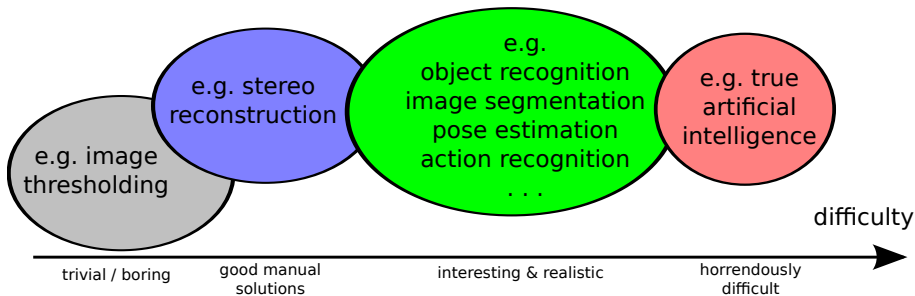
- For each example, correct output must be better than all others:

$$” \quad F(x^i, y^i) \geq \Delta(y^i, y) + F(x^i, y) \quad \text{for all } y \in \mathcal{Y} \setminus \{y^i\}. \quad ”$$

- convex optimization problem (similar to multiclass SVM)
- many equivalent formulations \rightarrow different training algorithms
- training calls MAP repeatedly, no probabilistic inference.**

For what problems can we learn solution?

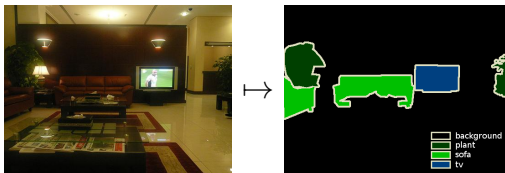
- we need training data,
- we need to be able to perform MAP/MAP-like prediction.



For many interesting problems, we cannot do exact MAP prediction.

The Need for Approximate Inference

- Image segmentation on pixel grid:
 - ▶ two labels and positive pairwise weights: exact MAP doable
 - ▶ more than two labels: MAP generally NP-hard



- Stereo reconstruction (grid) $\#labels = \#depth\ levels$:
 - ▶ convex penalization of disparities: exact MAP doable
 - ▶ non-convex (truncated) penalization: NP-hard
- Human Pose estimation
 - ▶ limbs independent given torso: exact MAP doable
 - ▶ limbs interact with each other: MAP NP-hard

Approximate inference $\hat{y} \approx \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y)$:

- has been used for decades, e.g. for
 - ▶ classical combinatorial problems: traveling salesman, knapsack
 - ▶ energy minimization: Markov/Conditional Random Fields
- often yield good (=almost perfect) solutions y
- can have guarantees $F(x, y^*) \leq F(x, \hat{y}) \leq (1 + \epsilon)F(x, y^*)$
- typically much faster than exact inference

Can't we use that for S-SVM training? **Iterating it is problematic!**

- In each S-SVM training iteration
 - ▶ we solve $y^* = \operatorname{argmax}_y [\Delta(y^i, y) + F(x^i, y)]$,
 - ▶ we use y^* to update w ,
 - ▶ we stop, when there are no more violated constraints / updates.
- With $\hat{y} \approx \operatorname{argmax}_y [\Delta(y^i, y) + F(x^i, y)]$
 - ▶ errors can accumulate,
 - ▶ we might stop too early.

Training S-SVMs with approximate MAP is active research question.

Summary – Learning for Structured Prediction

Task: Learn parameter w for a function $f : \mathcal{X} \rightarrow \mathcal{Y}$

$$f(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \langle w, \varphi(x, y) \rangle$$

Regularized risk minimization framework:

- loss function:
 - ▶ logistic loss \rightarrow conditional random fields (CRF)
 - ▶ Hinge loss \rightarrow structured support vector machine (S-SVM)
- regularization:
 - ▶ L^2 norm of $w \leftrightarrow$ Gaussian prior

Neither guarantees better prediction accuracy than the other.

see e.g. [Nowozin et al., ECCV 2010]

Difference is in numeric optimization:

- Use CRFs if you can do probabilistic inference.
- Use S-SVMs if you can do loss-augmented MAP prediction.

Open Research in Structured Learning

- Faster training
 - ▶ CRFs need many runs of probabilistic inference,
 - ▶ SSVMs need many runs of MAP-like predictions.
- Reduce amount of training data necessary
 - ▶ semi-supervised learning? transfer learning?
- Understand competing training methods
 - ▶ when to use *probabilistic training*, when *maximum margin*?
- **Understand S-SVM with approximate inference**
 - ▶ very often, exactly solving $\operatorname{argmax}_y F(x, y)$ is infeasible.
 - ▶ can we guarantee convergence even with approximate inference?

Vacancies at I.S.T. Austria, Vienna



More info: www.ist.ac.at
or chl@ist.ac.at.

- **PhD at I.S.T. Graduate School**
 - 1(2) + 3 yr PhD program
 - full scholarship
 - flexible starting dates
- **PostDoc in my group**
 - *computer vision*
 - ▶ object/attribute prediction.
 - *machine learning*
 - ▶ structured output learning.
 - curiosity driven basic research
 - ▶ no project deadlines,
 - ▶ no mandatory teaching, ...
- **Professor / Assistant Professor**
- **Visiting Scientist**
- **Internships**