

Algorithms for Predicting Structured Data

Thomas Gärtner / Shankar Vembu

Fraunhofer IAIS / UIUC

ECML PKDD 2010

Structured Prediction

Predicting multiple outputs with complex internal structure and dependencies

Contrast with single-valued prediction problems like binary classification and regression

Part-of-Speech Tagging

Input: Sentence

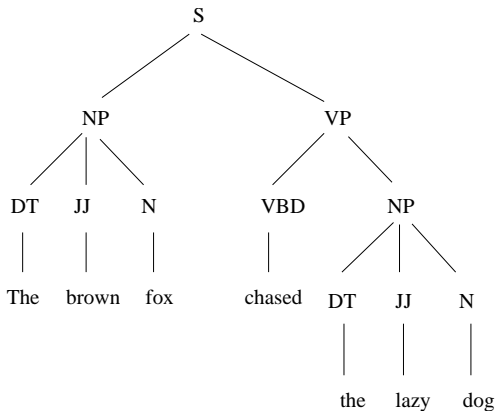
Output: Part-of-speech tags

The	brown	fox	chased	the	lazy	dog
DT	JJ	N	VBD	DT	JJ	N

Linguistic Parsing

Input: Sentence

Output: Parse tree



Machine Learning Problems

Multi-class prediction

Multi-label classification

Hierarchical classification

Label ranking

...

Real World Applications

Natural language processing

Computational biology

Bioinformatics

Computer vision

Robotics

...

Outline

Part I: From Binary to Structured Prediction

- ▶ Loss Functions
- ▶ Algorithms

Part II: Exact and Approximate Inference

- ▶ Approximate Inference: + / - ve Results
- ▶ Complexity of Learning
- ▶ New Training Methods

Part III: Weakly Supervised Learning

Tutorial Focus

Discriminative methods in structured prediction

Algorithmic techniques

Topics Not Covered

Advanced optimization methods

Connections to reinforcement learning

Learning theory / Generalization bounds

Inference in graphical models

X-supervised learning

Outline

Part I: Binary to Structured
Loss Functions
Algorithms

Road Map

Perceptron → Structured perceptron

Regression → Kernel dependency estimation

Support Vector Machines → Structured SVMs

Logistic Regression → Conditional random fields

Preliminaries

Supervised Learning

Input and output spaces \mathcal{X}, \mathcal{Y}

Classification $\mathcal{Y} = \{-1, +1\}$; Regression $\mathcal{Y} = \mathbb{R}$

Training data $X = (\mathbf{x}_1, \dots, \mathbf{x}_m) \in \mathcal{X}^m$
 $Y = (y_1, \dots, y_m) \in \mathcal{Y}^m$

(drawn i.i.d. from an unknown distribution)

Function approximation $f : \mathcal{X} \rightarrow \mathcal{Y}$

Good performance on unseen examples where performance is measured w.r.t. a loss function

Generative and Discriminative Learning

Generative learning (e.g., Naïve Bayes, HMM)

- ▶ Model joint probability $p(x, y)$
- ▶ Predict using Bayes rule

$$p(y | x) = \frac{p(x, y)}{\sum_{z \in \mathcal{Y}} p(x, z)} = \frac{p(x | y)p(y)}{\sum_{z \in \mathcal{Y}} p(x, z)}$$

Discriminative learning (e.g., Logistic regression, CRF)

- ▶ Model conditional probability $p(y | x)$ OR
- ▶ Mapping from inputs to outputs $f : \mathcal{X} \rightarrow \mathcal{Y}$

Regularized Risk Minimization

$$\operatorname{argmin}_{f \in \mathcal{H}} \sum_{i=1}^m \ell(f(x_i), y_i) + \lambda \Omega(f)$$

$\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$, non-negative loss function

$\Omega : \mathcal{H} \rightarrow \mathbb{R}_+$, regularization function

$\lambda > 0$, regularization parameter

Regularized Risk Minimization

Example: Linear classifiers $f(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle$

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \sum_{i=1}^m \ell(\langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle, y_i) + \lambda \|\mathbf{w}\|_2^2$$

$\phi : \mathcal{X} \rightarrow \mathbb{R}^n$, feature map

Structured Prediction

Input and output spaces \mathcal{X}, \mathcal{Y}

Training data $X = (\mathbf{x}_1, \dots, \mathbf{x}_m) \in \mathcal{X}^m$
 $Y = (\mathbf{y}_1, \dots, \mathbf{y}_m) \in \mathcal{Y}^m$

Joint scoring function $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$

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

Joint Feature Maps

Extract features from inputs AND outputs

Feature map, $\phi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{H}$

Joint input-output kernel,

$$k[(x, y), (x', y')] = \langle \phi(x, y), \phi(x', y') \rangle$$

Joint Feature Maps

Example: Multi-label classification, Hierarchical classification

$$\mathbf{x} \in \mathbb{R}^n, \quad \mathbf{y} \in \{0, 1\}^d$$

$$\phi(\mathbf{x}, \mathbf{y}) = \mathbf{x} \otimes \mathbf{y}$$

Tensor product kernel: $k[(\mathbf{x}, \mathbf{y}), (\mathbf{x}', \mathbf{y}')] = k_{\mathcal{X}}(\mathbf{x}, \mathbf{x}')k_{\mathcal{Y}}(\mathbf{y}, \mathbf{y}')$

Loss Functions

(Binary \rightarrow Structured)

Regularized Risk Minimization

$$\operatorname{argmin}_{f \in \mathcal{H}} \sum_{i=1}^m \ell(f(x_i), y_i) + \lambda \Omega(f)$$

$\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$, non-negative loss function

$\Omega : \mathcal{H} \rightarrow \mathbb{R}_+$, regularization function

$\lambda > 0$, regularization parameter

Zero-One Loss

Binary:

$$\ell_{0-1}(y, z) = \begin{cases} 0 & \text{if } y = z \\ 1 & \text{otherwise} \end{cases}$$

Structured:

$$\ell_{\max}^{\Delta}(f, (x, y)) = \Delta(\operatorname{argmax}_{z \in \mathcal{Y}} f(x, z), y)$$

Non-convex, non-differentiable, hard to optimize

Use surrogate losses instead; convex upper bounds

Squared Loss

$$\ell_{\text{square}}(\mathbf{y}, \mathbf{z}) = (\mathbf{y} - \mathbf{z})^2, \quad \mathcal{Y} = \mathbb{R}$$

Used in regression problems

Extension to structured prediction is non-trivial due to inter-dependencies among the multiple outputs (more on this later)

Hinge Loss

Binary:

$$\ell_{\text{hinge}}(y, z) = \max(0, 1 - yz)$$

Structured:

$$\ell_{\text{hinge}}^{\Delta}(f, (x, y)) = \max_{z \in \mathcal{Y}} [\Delta(z, y) + f(x, z) - f(x, y)]$$

Logistic Loss

Binary:

$$\ell_{\log}(y, z) = \ln(1 + \exp(-yz))$$

Structured:

$$\ell_{\log}(f, (x, y)) = \ln \left[\sum_{z \in \mathcal{Y}} \exp(f(x, z)) \right] - f(x, y)$$

Exponential Loss

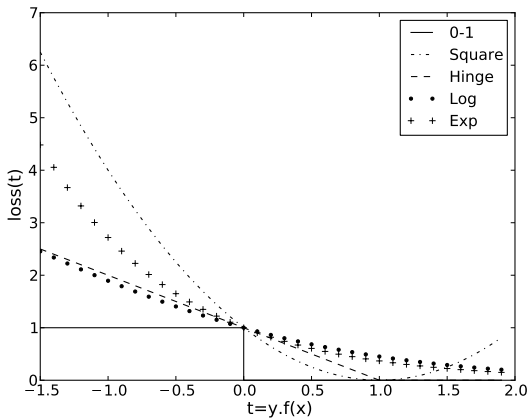
Binary:

$$l_{\text{exp}}(y, z) = \exp(-yz)$$

Structured:

$$l_{\text{exp}}(f, (x, y)) = \sum_{z \in \mathcal{Y}} \exp [f(x, z) - f(x, y)]$$

Loss Functions



Learning Algorithms

(parameter estimation)

Perceptron → Structured perceptron

Perceptron

Online learning algorithm

Prediction $\hat{y} = f(\mathbf{x}) = \text{sgn} \langle \mathbf{w}, \mathbf{x} \rangle$

Algorithm:

Initialize $\mathbf{w} = \mathbf{0}$

For $t = 1 \dots T$

1. receive input \mathbf{x}_t
2. predict $\hat{y}_t = \text{sgn} \langle \mathbf{w}, \mathbf{x}_t \rangle$
3. receive true label $y_t \in \{-1, +1\}$
4. if $y_t \neq \hat{y}_t$, then $\mathbf{w} \leftarrow \mathbf{w} + y_t \cdot \mathbf{x}_t$

Mistake Bound

Separable data [Block, 62; Novikoff, 62]

Theorem

Let $(x_1, y_1), \dots, (x_m, y_m)$ be a sequence of training examples with $\|x_i\| \leq R$ for all $i \in \llbracket m \rrbracket$. Suppose there exists a unit norm vector u such that $y_i(\langle u, x_i \rangle) \geq \gamma$ for all the examples. Then the number of mistakes made by the perceptron algorithm on this sequence is at most $(R/\gamma)^2$.

Mistake Bound

Inseparable data [Freund and Schapire, 99]

Theorem

Let $(x_1, y_1), \dots, (x_m, y_m)$ be a sequence of training examples with $\|x_i\| \leq R$ for all $i \in \llbracket m \rrbracket$. Let u be any unit norm weight vector and let $\gamma > 0$. Define the deviation of each example as $d_i = \max[0, \gamma - y_i(\langle u, x_i \rangle)]$ and define $D = \sqrt{\sum_{i=1}^m d_i^2}$. Then the number of mistakes made by the perceptron algorithm on this sequence of examples is at most

$$\frac{(R + D)^2}{\gamma^2} .$$

Constraint Classification

(Har-Peled et al., ALT 2002)

Generalizes perceptron for various problems including multiclass prediction, multilabel classification and label ranking

Example: Multiclass prediction, $\mathcal{Y} = \{1, \dots, d\}$

Weights $(w_1, \dots, w_d) \in \mathbb{R}^{nd}$

Prediction $\hat{y} = f(x) = \operatorname{argmax}_{i \in [d]} \langle w_i, x \rangle$

Constraint Classification

Algorithm:

Initialize $(w_1, \dots, w_d) = \mathbf{0}$

For $t = 1 \dots T$

1. receive input x_t and label $y_t \in \llbracket d \rrbracket$
2. construct $\tilde{y}_t = \{(y_t, i) \mid i \in \llbracket d \rrbracket \setminus y_t\}$
3. for all $(p, q) \in \tilde{y}_t$,
if $\langle w_p - w_q, x_t \rangle < 0$, then
 $w_p \leftarrow w_p + x$,
 $w_q \leftarrow w_q - x$

Structured Perceptron

(Collins, EMNLP 2002)

Linear **scoring function** $f(\mathbf{x}, y) = \langle \mathbf{w}, \phi(\mathbf{x}, y) \rangle$

Algorithm:

Initialize $\mathbf{w} = \mathbf{0}$

For $t = 1 \dots T$

1. receive input \mathbf{x}_t
2. predict $\hat{y}_t = \operatorname{argmax}_{y \in \mathcal{Y}} f(\mathbf{x}_t, y)$
3. receive true label $y_t \in \mathcal{Y}$
4. $\mathbf{w} \leftarrow \mathbf{w} + \phi(\mathbf{x}_t, y_t) - \phi(\mathbf{x}_t, \hat{y}_t)$

Mistake Bound

Separable data [Collins, 02]

Theorem

Let $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$ be a sequence of training examples which is separable with margin γ . Let R denote a constant that satisfies $\|\phi(x_i, y_i) - \phi(x_i, z)\| \leq R$, for all $i \in \llbracket m \rrbracket$, and for all $z \in \overline{GEN}(x_i)$. Then the number of mistakes made by the perceptron algorithm on this sequence of examples is at most R^2/γ^2 .

Mistake Bound

Inseparable data [Collins, 02]

For a training example (x_i, y_i) and a v, γ pair, define

$m_i = \langle v, \phi(x_i, y_i) \rangle - \max_{z \in \text{GEN}(x_i)} \langle v, \phi(x_i, z) \rangle$ and

$\epsilon_i = \max\{0, \gamma - m_i\}$ and define $D_{v, \gamma} = \sqrt{\sum_{i=1}^m \epsilon_i^2}$. Then the number of mistakes made by the structured perceptron is at most

$$\min_{v, \gamma} \frac{(R + D_{v, \gamma})^2}{\gamma^2} .$$

Regression \rightarrow KDE

Regularized Least-Squares Regression

$$\text{OPT: } \min_{w \in \mathbb{R}^n} \lambda \|w\|^2 + \sum_{i=1}^m (w^\top x_i - y_i)^2$$

$$\text{Solution: } w^* = (X^\top X + \lambda I)^{-1} X^\top y$$

where

$X \in \mathbb{R}^{m \times n}$: data/input matrix

$y \in \mathbb{R}^{m \times 1}$: label/output vector

I : identity matrix

Kernelized Version

$$\text{OPT1: } \min_{f \in \mathcal{H}} \lambda \|f\|^2 + \sum_{i=1}^m (f(x_i) - y_i)^2$$

$$\text{Representer theorem: } f^* = \sum_{i=1}^m c_i k(x_i, \cdot)$$

$$\text{OPT2: } \min_{\mathbf{c} \in \mathbb{R}^m} \lambda \mathbf{c}^\top \mathbf{K} \mathbf{c} + \|\mathbf{y} - \mathbf{K} \mathbf{c}\|^2$$

$$\text{Solution: } \mathbf{c}^* = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

Kernel Dependency Estimation

(Weston et al., NIPS 2002)

Output "feature" map $\psi : \mathcal{Y} \rightarrow \mathcal{H}_\mathcal{Y}$, $k_\mathcal{Y}(\mathbf{y}, \mathbf{y}') = \langle \psi(\mathbf{y}), \psi(\mathbf{y}') \rangle$

Note: Possible to consider a large class of non-linear loss functions in the output space

KDE uses kernel PCA to "decorrelate" the outputs, and trains independent RLSRs

Kernel Dependency Estimation

(Weston et al., NIPS 2002)

Algorithm:

1. decompose $\{\psi(y_i) \mid i \in \llbracket m \rrbracket\}$ into k orthogonal directions using KPCA
2. learn $f_j : \phi \rightarrow \mathbb{R}_j$ for each direction $j \in \llbracket k \rrbracket$ using RLSR
3. predict by solving the **pre-image problem**:

$$\hat{y}(x) = \operatorname{argmin}_{y \in \mathcal{Y}} \left([v_1^\top \psi(y), \dots, v_k^\top \psi(y)] - [f_1(x), \dots, f_k(x)] \right)^2$$

Kernel Dependency Estimation

(Weston et al., NIPS 2002)

Note: Projection of $\psi(y)$ onto the p^{th} principal component $v_p = \sum_{i=1}^m \alpha_i^p \psi(y_i)$ is given by $v_p^\top \psi(y) = \sum_{i=1}^m \alpha_i^p k_y(y_i, y)$, where α^p is the p^{th} eigenvector of the kernel matrix K_y (after centering it)

SVM \rightarrow Structured SVM

Support Vector Machine

Linear **large-margin** classifier $f(x) = \langle w, \phi(x) \rangle$

Minimizes the **hinge loss**

OPT:
$$\min_w \lambda \|w\|^2 + \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y_i \langle w, \phi(x_i) \rangle\}$$

OPT with slack variables:

$$\begin{aligned} \min_{w, \xi} \quad & \lambda \|w\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \\ \text{s.t. :} \quad & y_i \langle w, \phi(x_i) \rangle \geq 1 - \xi_i, \quad \forall i \in \llbracket m \rrbracket \\ & \xi_i \geq 0, \quad \forall i \in \llbracket m \rrbracket \end{aligned}$$

Support Vector Machine

SVM dual

$$\begin{aligned} \min_{\mathbf{c} \in \mathbb{R}^m} \quad & \frac{1}{2} \mathbf{c}^\top \mathbf{Y} \mathbf{K} \mathbf{Y} \mathbf{c} - \mathbf{1}^\top \mathbf{c} \\ \text{s.t.} \quad & 0 \leq c_i \leq \frac{1}{m\lambda}, \quad \forall i \in \llbracket m \rrbracket \end{aligned}$$

Representer theorem: $f(\cdot) = \sum_{i=1}^m c_i k(\mathbf{x}_i, \cdot)$

Structured SVM / Max-margin Markov Network

Minimizes **structured hinge loss**

Two formulations:

1. Slack rescaling
2. Margin rescaling

Structured SVM

Slack rescaling

$$\min_{\mathbf{w}, \xi} \quad \lambda \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i$$

$$\text{s.t. :} \quad \langle \mathbf{w}, \phi(\mathbf{x}_i, \mathbf{y}_i) \rangle - \langle \mathbf{w}, \phi(\mathbf{x}_i, \mathbf{z}) \rangle \geq 1 - \frac{\xi_i}{\Delta(\mathbf{y}_i, \mathbf{z})}, \quad \forall \mathbf{z}, \forall i$$

$$\xi_i \geq 0, \quad \forall i \in \llbracket m \rrbracket$$

Structured SVM

Margin rescaling

$$\min_{\mathbf{w}, \xi} \quad \lambda \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i$$

$$\text{s.t. :} \quad \langle \mathbf{w}, \phi(\mathbf{x}_i, \mathbf{y}_i) \rangle - \langle \mathbf{w}, \phi(\mathbf{x}_i, \mathbf{z}) \rangle \geq \Delta(\mathbf{y}_i, \mathbf{z}) - \xi_i, \quad \forall \mathbf{z}, \forall i$$

$$\xi_i \geq 0, \quad \forall i \in \llbracket m \rrbracket$$

Solving the OPT

Major issue: **Exponential** number of constraints

Suffices to design a sub-routine (loss-augmented inference) to compute

$$\hat{y} = \operatorname{argmax}_{z \in \mathcal{Y}} [1 - \langle w, \phi(x, y) - \phi(x, z) \rangle] \Delta(z, y)$$

$$\hat{y} = \operatorname{argmax}_{z \in \mathcal{Y}} [\Delta(z, y) - \langle w, \phi(x, y) - \phi(x, z) \rangle]$$

Iteratively add constraints

Cutting-Plane Method

(Tsochantaridis et al., JMLR 2005)

- 1: Input: $(x_1, y_1), \dots, (x_m, y_m), \lambda, \epsilon$
- 2: $S_i \leftarrow \emptyset$, for all $i \in \llbracket m \rrbracket$
- 3: **repeat**
- 4: **for** $i = 1, \dots, m$ **do**
- 5: $H(y) \equiv \Delta(y_i, y) + w^\top \phi(x_i, y) - w^\top \phi(x_i, y_i)$
- 6: **compute** $\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} H(y)$
- 7: **compute** $\xi_i = \max\{0, \max_{y \in S_i} H(y)\}$
- 8: **if** $H(\hat{y}) > \xi_i + \epsilon$ **then**
- 9: $S_i \leftarrow S_i \cup \{\hat{y}\}$
- 10: $w \leftarrow$ optimize primal over $\bigcup_i S_i$
- 11: **end if**
- 12: **end for**
- 13: **until** no S_i has changed during iteration

Theoretical Guarantees with Exact Inference

(Tsochantaridis et al., 2005; Finley and Joachims, ICML 2008)

Polynomial time termination Algorithm terminates in a polynomial number of iterations

Correctness Algorithm solves OPT accurately to a desired precision ϵ

Empirical risk bound $\frac{1}{m} \sum_{i=1}^m \xi_i$ upper bounds empirical risk

Kernelized SSVM

$$\min_{\alpha \in \mathbb{R}^{m|\mathcal{Y}|}} \sum_{i,j \in \llbracket m \rrbracket, z, z' \in \mathcal{Y}} \alpha_{iz} \alpha_{jz'} k[(x_i, z), (x_j, z')] - \sum_{i,z} \Delta(y_i, z) \alpha_{iz}$$

$$\text{s.t. : } \sum_{z \in \mathcal{Y}} \alpha_{iz} \leq \lambda, \forall i \in \llbracket m \rrbracket$$
$$\alpha_{iz} \geq 0, \forall i \in \llbracket m \rrbracket, \forall z \in \mathcal{Y}$$

Representer theorem:

$$f(\cdot, \cdot) = \sum_{i \in \llbracket m \rrbracket, z \in \mathcal{Y}} \alpha_{iz} k[(x_i, z), (\cdot, \cdot)]$$

Structured Perceptron Revisited

Linear scoring function $f(\mathbf{x}, y) = \langle \mathbf{w}, \phi(\mathbf{x}, y) \rangle$

Algorithm:

Initialize $\mathbf{w} = \mathbf{0}$

For $t = 1 \dots T$

1. receive input \mathbf{x}_t
2. predict $\hat{y}_t = \operatorname{argmax}_{y \in \mathcal{Y}} f(\mathbf{x}_t, y)$
3. receive true label $y_t \in \mathcal{Y}$
4. $\mathbf{w} \leftarrow \mathbf{w} + \phi(\mathbf{x}_t, y_t) - \phi(\mathbf{x}_t, \hat{y}_t)$

Note: Update rule is not influenced by the loss function $\Delta(y, z)$

Structured Perceptron Revisited

Replace **inference** with **loss-augmented inference**

Algorithm:

Initialize $w = \mathbf{0}$

For $t = 1 \dots T$

1. receive example pair \mathbf{x}_t, y_t
2. predict $\hat{y}_t = \operatorname{argmax}_{y \in \mathcal{Y}} [f(\mathbf{x}_t, y) + \Delta(y_t, y)]$
3. $w \leftarrow w + \eta_t (\phi(\mathbf{x}_t, y_t) - \phi(\mathbf{x}_t, \hat{y}_t))$

Min-Max Formulation

(Taskar et al., ICML 2005)

Recall brute force enumeration

$$\min_{\mathbf{w}, \xi} \quad \lambda \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i$$

$$\text{s.t. :} \quad \langle \mathbf{w}, \phi(\mathbf{x}_i, \mathbf{y}_i) \rangle - \langle \mathbf{w}, \phi(\mathbf{x}_i, \mathbf{z}) \rangle \geq \Delta(\mathbf{y}_i, \mathbf{z}) - \xi_i, \quad \forall \mathbf{z}, \forall i$$

$$\xi_i \geq 0, \quad \forall i \in \llbracket m \rrbracket$$

Min-Max Formulation

(Taskar et al., ICML 2005)

$$\min_{\mathbf{w}, \xi} \quad \lambda \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i$$

$$\text{s.t. :} \quad \langle \mathbf{w}, \phi(\mathbf{x}_i, \mathbf{y}_i) \rangle + \xi_i \geq \max_{\mathbf{z} \in \mathcal{Y}} [\langle \mathbf{w}, \phi(\mathbf{x}_i, \mathbf{z}) \rangle + \Delta(\mathbf{y}_i, \mathbf{z})], \quad \forall i \in \llbracket m \rrbracket$$

$$\xi_i \geq 0, \quad \forall i \in \llbracket m \rrbracket$$

Key steps:

- ▶ Plug-in LP inference
- ▶ Use LP duality to replace max with min
- ▶ Rewrite the OPT as a *concise* QP with polynomial number of variables and constraints (depends on the output structure)

Logistic Regression \rightarrow CRF

Logistic Regression

Probabilistic (binary) classifier

Likelihood function:

$$p(y | x, w) = \frac{\exp(y \langle \phi(x), w \rangle)}{\exp(y \langle \phi(x), w \rangle) + \exp(-y \langle \phi(x), w \rangle)}$$

Estimate parameters w by minimizing negative log-likelihood

Exponential Family

Family of probability distributions

$$p(\mathbf{x}; \mathbf{w}) = \exp(\langle \phi(\mathbf{x}), \mathbf{w} \rangle - \ln Z(\mathbf{w}))$$

$\phi(\mathbf{x})$ - Sufficient statistics of x

$Z(\mathbf{w})$ - Partition function

$$Z(\mathbf{w}) = \int_{\mathbf{x}} \exp(\langle \phi(\mathbf{x}), \mathbf{w} \rangle) d\mathbf{x}$$

Log-Partition Function

$$g(w) = \ln Z(w) = \ln \int_{\mathbf{x}} \exp(\langle \phi(\mathbf{x}), w \rangle) d\mathbf{x}$$

$g(w)$ is a cumulant generator, i.e.,

$$\partial_w g(w) = \frac{\int \phi(\mathbf{x}) \exp(\langle \phi(\mathbf{x}), w \rangle) d\mathbf{x}}{\int \exp(\langle \phi(\mathbf{x}), w \rangle) d\mathbf{x}} = \mathbf{E}_{\mathbf{x} \sim p(\mathbf{x}; w)} [\phi(\mathbf{x})]$$

$$\partial_w^2 g(w) = \mathbf{Cov}_{\mathbf{x} \sim p(\mathbf{x}; w)} [\phi(\mathbf{x})]$$

Note: $g(w)$ is convex!

Examples

Binomial

Multinomial

Gaussian

Laplace

Poisson

Dirichlet

...

(Conditional) Exponential Family

Family of conditional distributions

$$p(y | x, w) = \exp(\langle \phi(x, y), w \rangle - \ln Z(w | x))$$

$\phi(x, y)$ - Joint sufficient statistics of x and y

$Z(w | x)$ - Partition function

$$Z(w | x) = \int_{y \in \mathcal{Y}} \exp(\langle \phi(x, y), w \rangle)$$

MAP Estimation

$$p(w | Y, X) \approx p(w, Y | X) = p(Y | X, w)p(w)$$

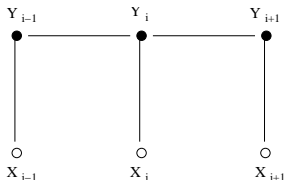
Point estimate with a Gaussian prior on w ,

$$p(w) = \exp(-\lambda \|w\|^2)$$

$$\begin{aligned}\hat{w} &= \operatorname{argmax}_w [\ln p(w, Y | X)] \\ &= \operatorname{argmax}_w \left[\frac{1}{m} \sum_{i=1}^m \ln p(y_i | x_i, w) - \lambda \|w\|^2 \right]\end{aligned}$$

Note: Convex optimization problem

Linear-Chain Conditional Random Fields

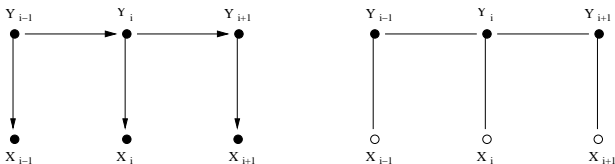


Cliques: (y_t, x_t) , (y_{t-1}, y_t) , for all t

$$p(y | x, w) = \exp \left(\sum_t [\langle \phi(y_t, x_t), w_{xy,t} \rangle + \langle \phi(y_{t-1}, y_t), w_{yy,t} \rangle] - \ln Z(w | x) \right)$$

- ▶ Efficient inference using dynamic programming
- ▶ MLE / MAP estimation of parameters

HMMs and CRFs



Hidden Markov Model: Maximize $p(x, y | w)$

Conditional Random Field: Maximize $p(y | x, w)$

Learning Reductions

(ICML 2009 tutorial)

Reductions: Transform complex learning problems into simpler, core problems

Desideratum: Good performance on the core problem should imply good performance on the complex problem

Examples

- ▶ Multi-class prediction → Binary classification
- ▶ Ranking → Classification
- ▶ Cost-sensitive classification → Binary classification
- ▶ Structured prediction → Binary classification (SEARN)

Structured \rightarrow Binary (SEARN)

(Daumé et al., 2009)

Decomposable outputs $y = (y_1, \dots, y_T)$

SEARN learns a policy π that maps tuples (x, y_1, \dots, y_{t-1}) to y_T

Reduces structured prediction to cost-sensitive (multi-class) classification

Good performance on cost-sensitive classification implies good performance on the structured prediction problem

Note: No **argmax** problem!

Reduction to Cost-Sensitive Classification

Distribution over cost-sensitive examples (input, c_1, \dots, c_K)

For every structured example (x, y)

1. sample t uniformly from $\llbracket T \rrbracket$
2. run policy π for $t - 1$ steps to yield $(\hat{y}_1, \dots, \hat{y}_{t-1})$
3. Input: $(x, \hat{y}_1, \dots, \hat{y}_{t-1})$
4. Costs:

$$c_k = \mathbf{E}_{\hat{y}_{t+1}, \dots, \hat{y}_T \sim \pi} \ell(y, (\hat{y}_1, \dots, \hat{y}_{t-1}, k, \hat{y}_{t-1}, \dots, \hat{y}_T)), \forall k \in \llbracket K \rrbracket$$

Intermediate Summary

Loss Functions (binary \rightarrow structured)

Discriminative Structured Prediction Algorithms

- ▶ Perceptron \rightarrow Structured perceptron
- ▶ RLSR \rightarrow KDE
- ▶ SVM \rightarrow Struct-SVM, M3N
- ▶ LogReg \rightarrow CRF
- ▶ Structured prediction \rightarrow Cost-sensitive classification