

Probabilistic Graphical Models in Computer Vision (IN2329)

Csaba Domokos

Summer Semester 2015/2016

11. Parameter learning	2
Loss function	3
Loss function	4
0/1 loss *	5
Hamming-loss *	6
Agenda for today's lecture *	7
Probabilistic parameter learning	8
Recall: Regularized maximum conditional likelihood training *	9
Numerical solution	10
Stochastic gradient descent	11
Stochastic gradient descent *	12
Using of the output structure	13
Two-stage learning	14
Piecewise learning	15
Piecewise learning	16
Loss-minimizing parameter learning	17
Loss-minimizing parameter learning	18

Regularized loss minimization 19

Digression: Support Vector Machine * 20

Digression: Support Vector Machine * 21

Redefining the loss function 22

Structured hinge loss 23

Structured Support Vector Machine 24

S-SVM: Toy example * 25

Subgradient * 26

Subgradient descent minimization * 27

Numerical solution 28

Calculating the subgradient. 29

Subgradient descent S-SVM learning * 30

Stochastic subgradient descent S-SVM learning * 31

Summary of S-SVM learning * 32

Next lecture: Summary of the course * 33

Literature * 34

Loss function

The goal is to make prediction $\mathbf{y} \in \mathcal{Y}$, *as good as possible*, about unobserved properties (e.g., class label) for a given data instance $\mathbf{x} \in \mathcal{X}$.

In order to measure quality of **prediction** $f : \mathcal{X} \rightarrow \mathcal{Y}$ we define a **loss function**

$$\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+,$$

that is $\Delta(\mathbf{y}, \mathbf{y}')$ measures the cost of predicting \mathbf{y}' when the correct label is \mathbf{y} .

Let us denote the *model distribution* by $p(\mathbf{y} | \mathbf{x})$ and the *true (conditional) data distribution* by $d(\mathbf{y} | \mathbf{x})$. The quality of prediction can be expressed by the **expected loss**:

$$\begin{aligned} \mathcal{R}_f^\Delta(\mathbf{x}) &:= \mathbb{E}_{\mathbf{y} \sim d(\mathbf{y} | \mathbf{x})} [\Delta(\mathbf{y}, f(\mathbf{x}))] = \sum_{\mathbf{y} \in \mathcal{Y}} d(\mathbf{y} | \mathbf{x}) \Delta(\mathbf{y}, f(\mathbf{x})) \\ &\approx \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y} | \mathbf{x})} [\Delta(\mathbf{y}, f(\mathbf{x}))], \end{aligned}$$

assuming that $p(\mathbf{y} | \mathbf{x}) \approx d(\mathbf{y} | \mathbf{x})$.

0/1 loss *

In general, the *loss function* is application dependent. Arguably one of the most common *loss functions* for labelling tasks is the **0/1 loss**, that is

$$\Delta_{0/1}(\mathbf{y}, \mathbf{y}') = \llbracket \mathbf{y} \neq \mathbf{y}' \rrbracket = \begin{cases} 0, & \text{if } \mathbf{y} = \mathbf{y}' \\ 1, & \text{otherwise.} \end{cases}$$

Minimizing the *expected loss* of the 0/1 loss yields

$$\begin{aligned} \mathbf{y}^* &\in \operatorname{argmin}_{\mathbf{y}' \in \mathcal{Y}} \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x})} [\Delta_{0/1}(\mathbf{y}, \mathbf{y}')] = \operatorname{argmin}_{\mathbf{y}' \in \mathcal{Y}} \sum_{\mathbf{y} \in \mathcal{Y}} p(\mathbf{y} | \mathbf{x}) \Delta_{0/1}(\mathbf{y}, \mathbf{y}') \\ &= \operatorname{argmin}_{\mathbf{y}' \in \mathcal{Y}} \sum_{\mathbf{y} \in \mathcal{Y}, \mathbf{y} \neq \mathbf{y}'} p(\mathbf{y} | \mathbf{x}) = \operatorname{argmin}_{\mathbf{y}' \in \mathcal{Y}} (1 - p(\mathbf{y}' | \mathbf{x})) = \operatorname{argmax}_{\mathbf{y}' \in \mathcal{Y}} p(\mathbf{y}' | \mathbf{x}) \\ &= \operatorname{argmin}_{\mathbf{y}' \in \mathcal{Y}} E(\mathbf{y}'; \mathbf{x}) . \end{aligned}$$

This shows that the *optimal* prediction $f(\mathbf{x}) = \mathbf{y}^*$ in this case is given by **MAP inference**.

Hamming-loss *

Another popular choice of loss function is the **Hamming-loss**, which counts the percentage of mis-labeled variables:

$$\Delta_H(\mathbf{y}, \mathbf{y}') = \frac{1}{|\mathcal{Y}|} \sum_{i \in \mathcal{Y}} \mathbb{1}[y_i \neq y'_i].$$

For example, in *pixel-wise image segmentation*, the **Hamming-loss** is proportional to the number of mis-classified pixels, whereas the 0/1 **loss** assigns the same cost to every labeling that is not pixel-by-pixel identical to the correct one.

The *expected loss* of the *Hamming-loss* takes the form (see Exercise)

$$\mathcal{R}_f^H(\mathbf{x}) = 1 - \frac{1}{|\mathcal{Y}|} p(Y_i = f(\mathbf{x})_i \mid \mathbf{x}),$$

which is minimized by predicting with $f(\mathbf{x})_i = \operatorname{argmax}_{y_i \in \mathcal{Y}_i} p(Y_i = y_i \mid \mathbf{x})$.

To evaluate this prediction rule, we rely on **probabilistic inference**.

Agenda for today's lecture *

Probabilistic parameter learning is the task of estimating the parameter \mathbf{w} that minimizes the **expected dissimilarity** of a parameterized *model distribution* $p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$ and the (*unknown*) conditional *data distribution* $d(\mathbf{y} \mid \mathbf{x})$:

$$\text{KL}_{\text{tot}}(d \parallel p) = \sum_{\mathbf{x} \in \mathcal{X}} d(\mathbf{x}) \sum_{\mathbf{y} \in \mathcal{Y}} d(\mathbf{y} \mid \mathbf{x}) \log \frac{d(\mathbf{y} \mid \mathbf{x})}{p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})}.$$

Loss minimizing parameter learning is the task of finding the parameter \mathbf{w} such that the *expected prediction loss*

$$\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim d(\mathbf{x}, \mathbf{y})} [\Delta(\mathbf{y}, f(\mathbf{x}))]$$

is as small as possible, where $f(\mathbf{x}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$ is a *prediction function*, $d(\mathbf{x}, \mathbf{y})$ is the (*unknown*) *true data distribution*, and $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ is a *loss function*.

Recall: Regularized maximum conditional likelihood training *

Let $p(\mathbf{y} \mid \mathbf{x}, \mathbf{w}) = \frac{1}{Z(\mathbf{x}, \mathbf{w})} \exp(-\langle \mathbf{w}, \varphi(\mathbf{x}, \mathbf{y}) \rangle)$ be a *probability distribution parameterized by* $\mathbf{w} \in \mathbb{R}^D$, and let $\mathcal{D} = \{(\mathbf{x}^n, \mathbf{y}^n)\}_{n=1, \dots, N}$ be a set of *i.i.d. training examples*. For any $\lambda > 0$, **regularized maximum conditional likelihood** training chooses the parameter as

$$\begin{aligned} \mathbf{w}^* &\in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} L(\mathbf{w}) \\ &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} \lambda \|\mathbf{w}\|^2 + \sum_{n=1}^N \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}^n) \rangle + \sum_{n=1}^N \log Z(\mathbf{x}^n, \mathbf{w}) . \end{aligned}$$

Numerical solution

$$\nabla_{\mathbf{w}} L(\mathbf{w}) = 2\lambda \mathbf{w} + \sum_{n=1}^N (\varphi(\mathbf{x}^n, \mathbf{y}^n) - \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y} \mid \mathbf{x}^n, \mathbf{w})} [\varphi(\mathbf{x}^n, \mathbf{y})]) .$$

In a naïve way, the complexity of the *gradient computation* is $\mathcal{O}(K^{|\mathcal{V}|}ND)$.

$$\lambda \|\mathbf{w}\|^2 + \sum_{n=1}^N \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}^n) \rangle + \sum_{n=1}^N \log Z(\mathbf{x}^n, \mathbf{w}) .$$

In a naïve way, the complexity of the *line search* is $\mathcal{O}(K^{|\mathcal{V}|}ND)$ (for each evaluation of L), where

- N is the number of data samples,
- D is the dimension of weight vector,
- $K = \max_{i \in \mathcal{V}} |\mathcal{Y}_i|$ is the (maximal) number of possible labels of each output nodes.

Stochastic gradient descent

If the training set \mathcal{D} is too large, one can create a random subset $\mathcal{D}' \subset \mathcal{D}$ and estimate the gradient $\nabla_{\mathbf{w}} L(\mathbf{w})$ on \mathcal{D}' . In an extreme case, one may *randomly select* only **one** sample and calculate the gradient

$$\tilde{\nabla}_{\mathbf{w}}^{(\mathbf{x}^n, \mathbf{y}^n)} L(\mathbf{w}) = 2\lambda \mathbf{w} + \varphi(\mathbf{x}^n, \mathbf{y}^n) - \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}^n, \mathbf{w})}[\varphi(\mathbf{x}^n, \mathbf{y})].$$

This approach is called **stochastic gradient descent** (SGD).

Note that line search is not possible, therefore, we need for an extra parameter, referred to as **step-size** η_t for each iteration ($t = 1, \dots, T$).

Stochastic gradient descent *

Input: Step-sizes η_1, \dots, η_T for all the T iterations.

Output: The learned weight vector $\mathbf{w} \in \mathbb{R}^D$.

- 1: $\mathbf{w} \leftarrow \mathbf{0}$
- 2: **for** $t = 1, \dots, T$ **do**
- 3: $(\mathbf{x}^n, \mathbf{y}^n) \leftarrow$ a randomly chosen training example
- 4: $\mathbf{d} \leftarrow -\tilde{\nabla}_{\mathbf{w}}^{(\mathbf{x}^n, \mathbf{y}^n)} L(\mathbf{w})$
- 5: $\mathbf{w} \leftarrow \mathbf{w} + \eta_t \mathbf{d}$
- 6: **end for**
- 7: **return** \mathbf{w}

If the step-size is chosen correctly (e.g., $\eta_t = \frac{\eta}{t}$), then SGD converges to $\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} L(\mathbf{w})$. However, it needs more iterations than *gradient descent*, but each iteration is (much) faster.

Using of the output structure

Assume a set of factors \mathcal{F} in a factor graph model, such that the vector $\varphi(\mathbf{x}, \mathbf{y})$ decomposes as $\varphi(\mathbf{x}, \mathbf{y}) = [\varphi_F(\mathbf{x}_F, \mathbf{y}_F)]_{F \in \mathcal{F}}$. Thus

$$\begin{aligned}\mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, \mathbf{w})}[\varphi(\mathbf{x}, \mathbf{y})] &= [\mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, \mathbf{w})}[\varphi_F(\mathbf{x}_F, \mathbf{y}_F)]]_{F \in \mathcal{F}} \\ &= [\mathbb{E}_{\mathbf{y}_F \sim p(\mathbf{y}_F|\mathbf{x}_F, \mathbf{w}_F)}[\varphi_F(\mathbf{x}_F, \mathbf{y}_F)]]_{F \in \mathcal{F}},\end{aligned}$$

where

$$\mathbb{E}_{\mathbf{y}_F \sim p(\mathbf{y}_F|\mathbf{x}_F, \mathbf{w}_F)}[\varphi_F(\mathbf{x}_F, \mathbf{y}_F)] = \sum_{\mathbf{y}_F \in \mathcal{Y}_F} p(\mathbf{y}_F | \mathbf{x}_F, \mathbf{w}_F) \varphi_F(\mathbf{x}_F, \mathbf{y}_F).$$

Factor marginals $\mu_F = p(\mathbf{y}_F | \mathbf{x}_F, \mathbf{w}_F)$ are generally (much) easier to calculate than the complete conditional distribution $p(\mathbf{y} | \mathbf{x}, \mathbf{w})$.

They can be either computed *exactly* (e.g., by applying Belief propagation yielding complexity $\mathcal{O}(K^{|\mathcal{F}_{\max}|} |\mathcal{V}| ND)$, where $|\mathcal{F}_{\max}| = \max_{F \in \mathcal{F}} |N(F)|$ is the maximal factor size) or *approximated*.

Two-stage learning

The idea here is to split learning of energy functions into two steps:

1. learning of unary energies via *classifiers*, and
2. learning of their importance and the weighting factors of pairwise (and higher-order) energies.

$$E(\mathbf{y}; \mathbf{x}) = \sum_{i \in \mathcal{V}} w_i E_i(y_i; x_i) + \sum_{(i,j) \in \mathcal{E}'} w_{ij} E_{ij}(y_i, y_j).$$

As an advantage, it results in a *faster* learning method. However, if local classifiers for E_i perform badly, then CRF learning **cannot** fix it.

Piecewise learning

Assume a set of factors \mathcal{F} in a *factor graph model*, such that $\varphi(\mathbf{x}, \mathbf{y}) = [\varphi_F(\mathbf{x}_F, \mathbf{y}_F)]_{F \in \mathcal{F}}$.

We now **approximate** $p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$ by a distribution that is a product over the factors:

$$p_{\text{PW}}(\mathbf{y} \mid \mathbf{x}, \mathbf{w}) := \prod_{F \in \mathcal{F}} p_F(\mathbf{y}_F \mid \mathbf{x}_F, \mathbf{w}_F) = \prod_{F \in \mathcal{F}} \frac{\exp(-\langle \mathbf{w}_F, \varphi_F(\mathbf{x}_F, \mathbf{y}_F) \rangle)}{Z_F(\mathbf{x}_F, \mathbf{w}_F)} .$$

By minimizing the negative conditional log-likelihood function $L(\mathbf{w})$, we get

$$\begin{aligned} \mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} L(\mathbf{w}) &\approx \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} \lambda \|\mathbf{w}\|^2 - \sum_{n=1}^N \log \prod_{F \in \mathcal{F}} p_F(\mathbf{y}_F^n \mid \mathbf{x}_F^n, \mathbf{w}_F) \\ &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} \sum_{F \in \mathcal{F}} \lambda \|\mathbf{w}_F\|^2 + \sum_{n=1}^N \langle \mathbf{w}_F, \varphi_F(\mathbf{x}_F^n, \mathbf{y}_F^n) \rangle + \sum_{n=1}^N \log Z_F(\mathbf{x}_F^n, \mathbf{w}_F) . \end{aligned}$$

Piecewise learning

$$\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} \sum_{F \in \mathcal{F}} \lambda \|\mathbf{w}_F\|^2 + \sum_{n=1}^N \langle \mathbf{w}_F, \varphi_F(\mathbf{x}_F^n, \mathbf{y}_F^n) \rangle + \sum_{n=1}^N \log Z_F(\mathbf{x}_F^n, \mathbf{w}_F) .$$

Consequently, piecewise training chooses the parameters $\mathbf{w}^* = [\mathbf{w}_F^*]_{F \in \mathcal{F}}$ as

$$\mathbf{w}_F^* \in \operatorname{argmin}_{\mathbf{w}_F \in \mathbb{R}} \lambda \|\mathbf{w}_F\|^2 + \sum_{n=1}^N \langle \mathbf{w}_F, \varphi_F(\mathbf{x}_F^n, \mathbf{y}_F^n) \rangle + \sum_{n=1}^N \log Z_F(\mathbf{x}_F^n, \mathbf{w}_F) .$$

One can perform gradient-based training for each factor as long as the individual factors remain small.

Comparing $p_{\text{PW}}(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$ with the exact $p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$, we see that the exact $Z(\mathbf{w})$ does not factorize into a product of simpler terms, whereas its piecewise approximation $Z_{\text{PW}}(\mathbf{w})$ factorizes over the set of factors.

The simplification made by piece-wise training of CRFs resembles **two-stage learning**.

Loss-minimizing parameter learning

Let $\mathcal{D} = \{(\mathbf{x}^1, \mathbf{y}^1), \dots, (\mathbf{x}^N, \mathbf{y}^N)\} \subseteq \mathcal{X} \times \mathcal{Y}$ be *i.i.d.* samples from the (unknown) *true data distribution* $d(\mathbf{x}, \mathbf{y})$ and $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ be a *loss function*. The task is to find a weight vector \mathbf{w} that leads to **minimal expected loss**

$$\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim d(\mathbf{x}, \mathbf{y})} [\Delta(\mathbf{y}, f(\mathbf{x}))]$$

for a *prediction function* $f(\mathbf{x}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} g(\mathbf{x}, \mathbf{y}; \mathbf{w})$, where $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is an **auxiliary function** that is parameterized by $\mathbf{w} \in \mathbb{R}^D$.

Pros:

- We directly optimize for the *quantity of interest*, i.e. the expected loss.
- We do not need to compute the *partition function* Z .

Cons:

- There is no probabilistic reasoning to find \mathbf{w} .
- We need to know the *loss function* already at training time.

Regularized loss minimization

Let us define the *auxiliary function* as

$$g(\mathbf{x}, \mathbf{y}; \mathbf{w}) := -E(\mathbf{y}; \mathbf{x}, \mathbf{w}) = -\langle \mathbf{w}, \varphi(\mathbf{x}, \mathbf{y}) \rangle .$$

We aim to find the parameter \mathbf{w}^* that minimizes

$$\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim d(\mathbf{x}, \mathbf{y})} [\Delta(\mathbf{y}, f(\mathbf{x}))] = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim d(\mathbf{x}, \mathbf{y})} [\Delta(\mathbf{y}, \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} g(\mathbf{x}, \mathbf{y}; \mathbf{w}))] .$$

However, $d(\mathbf{x}, \mathbf{y})$ is unknown, hence we apply *approximation*:

$$\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim d(\mathbf{x}, \mathbf{y})} [\Delta(\mathbf{y}, \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} g(\mathbf{x}, \mathbf{y}; \mathbf{w}))] \approx \frac{1}{N} \sum_{n=1}^N \Delta(\mathbf{y}^n, \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} g(\mathbf{x}^n, \mathbf{y}^n; \mathbf{w})) .$$

Moreover, we add the **regularizer** $\lambda \|\mathbf{w}\|^2$ in order to avoid *overfitting*.

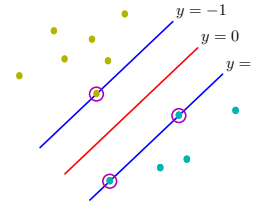
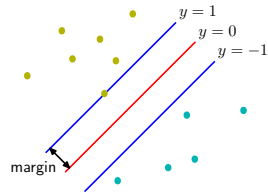
Therefore, we get a new objective, that is

$$\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} \lambda \|\mathbf{w}\|^2 + \frac{1}{N} \sum_{n=1}^N \Delta(\mathbf{y}^n, \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} g(\mathbf{x}^n, \mathbf{y}^n; \mathbf{w})) .$$

Digression: Support Vector Machine *

Let us consider the **binary classification** problem. Suppose we are given a set of labeled points $\{(\mathbf{x}^1, t^1), \dots, (\mathbf{x}^N, t^N)\}$, where $\mathbf{x}^n \in \mathbb{R}^D$ and $t^n \in \{-1, 1\}$ for all $n = 1, \dots, N$.

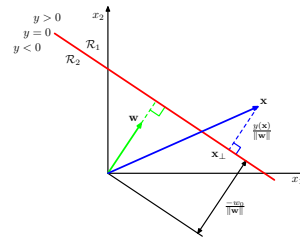
The *goal* is to find a hyperplane $y(\mathbf{x}) := \langle \mathbf{w}, \mathbf{x} \rangle + w_0$ separating the input data according to their label t^n .



More precisely, $y(\mathbf{x}^n) > 0$ for points having $t^n = 1$ and $y(\mathbf{x}^n) < 0$ for points having $t^n = -1$, that is $t^n \cdot y(\mathbf{x}^n) \geq 1$ for all training points.

If such a hyperplane exists, then we say the *training set* is **linearly separable**.

Digression: Support Vector Machine *



We want to solve the following minimization problem:

$$\mathbf{w}^* \in \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{w}\|, \text{ subject to } t^n(\langle \mathbf{w}, \mathbf{x}^n \rangle + w_0) \geq 1, \text{ for all } n = 1, \dots, N.$$

Since the training set is not necessarily *linearly separable*, instead, we consider the following minimization for $\lambda > 0$

$$\mathbf{w}^* \in \underset{\mathbf{w}}{\operatorname{argmin}} \lambda \|\mathbf{w}\| + \frac{1}{N} \sum_{n=1}^N \max(0, 1 - t^n(\langle \mathbf{w}, \mathbf{x}^n \rangle + w_0)).$$

where $\ell(y) = \max(0, 1 - ty)$ is called the **hinge loss** function.

Redefining the loss function

$$\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} \lambda \|\mathbf{w}\|^2 + \frac{1}{N} \sum_{n=1}^N \Delta(\mathbf{y}^n, \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} g(\mathbf{x}^n, \mathbf{y}^n; \mathbf{w})) .$$

Note that the loss function $\Delta(\mathbf{y}, \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} g(\mathbf{x}, \mathbf{y}; \mathbf{w}))$ is piecewise constant, hence it is **discontinuous**, therefore we cannot use gradient-based techniques. As a remedy we will *replace* $\Delta(\mathbf{y}, \mathbf{y}')$ with well behaved function $\ell(\mathbf{x}, \mathbf{y}; \mathbf{w})$, i.e. it is *continuous* and *convex* with respect to w .

Typically, ℓ is chosen such that it is an **upper bound** to Δ .

Therefore, we get a new objective, that is

$$\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} \lambda \|\mathbf{w}\|^2 + \frac{1}{N} \sum_{n=1}^N \ell(\mathbf{x}^n, \mathbf{y}^n; \mathbf{w}) .$$

Structured hinge loss

Let $\bar{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} g(\mathbf{x}^n, \mathbf{y}; \mathbf{w})$, then

$$\begin{aligned} \Delta(\mathbf{y}^n, \bar{\mathbf{y}}) &\leq \Delta(\mathbf{y}^n, \bar{\mathbf{y}}) + g(\mathbf{x}^n, \bar{\mathbf{y}}; \mathbf{w}) - g(\mathbf{x}^n, \mathbf{y}^n; \mathbf{w}) \\ &\leq \max_{\mathbf{y} \in \mathcal{Y}} (\Delta(\mathbf{y}^n, \mathbf{y}) + g(\mathbf{x}^n, \mathbf{y}; \mathbf{w}) - g(\mathbf{x}^n, \mathbf{y}^n; \mathbf{w})) \\ &\triangleq \ell(\mathbf{x}^n, \mathbf{y}^n, \mathbf{w}) . \end{aligned}$$

The **structured hinge loss** ℓ provides an upper bound for the *loss function* Δ . Note that ℓ is continuous and convex, since it is a maximum over *affine functions*.

We remark that

$$\begin{aligned} \ell(\mathbf{x}^n, \mathbf{y}^n, \mathbf{w}) &\triangleq \max_{\mathbf{y} \in \mathcal{Y}} (\Delta(\mathbf{y}^n, \mathbf{y}) + g(\mathbf{x}^n, \mathbf{y}; \mathbf{w}) - g(\mathbf{x}^n, \mathbf{y}^n; \mathbf{w})) \\ &= \max \left(0, \max_{\mathbf{y} \in \mathcal{Y}} (\Delta(\mathbf{y}^n, \mathbf{y}) + g(\mathbf{x}^n, \mathbf{y}; \mathbf{w}) - g(\mathbf{x}^n, \mathbf{y}^n; \mathbf{w})) \right) \\ &= \max_{\mathbf{y} \in \mathcal{Y}} (\Delta(\mathbf{y}^n, \mathbf{y}) - \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}) \rangle + \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}^n) \rangle) . \end{aligned}$$

Structured Support Vector Machine

Let $g(\mathbf{x}, \mathbf{y}; \mathbf{w}) = -\langle \mathbf{w}, \varphi(\mathbf{x}, \mathbf{y}) \rangle$ be an *auxiliary function* parameterized by $\mathbf{w} \in \mathbb{R}^D$. For any $C > 0$, **structured support vector machine** (S-SVM) training chooses the parameter

$$\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{N} \sum_{n=1}^N \ell(\mathbf{x}^n, \mathbf{y}^n, \mathbf{w})$$

with

$$\ell(\mathbf{x}^n, \mathbf{y}^n, \mathbf{w}) = \max_{\mathbf{y} \in \mathcal{Y}} (\Delta(\mathbf{y}^n, \mathbf{y}) - \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}) \rangle + \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}^n) \rangle) .$$

Both probabilistic parameter learning and S-SVM do *regularized risk minimization*. For probabilistic parameter learning, the *regularized conditional log-likelihood function* can be written as:

$$\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} \frac{\|\mathbf{w}\|^2}{2\sigma^2} + \sum_{n=1}^N \log \sum_{\mathbf{y} \in \mathcal{Y}} \exp (\langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}) \rangle - \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}^n) \rangle) .$$

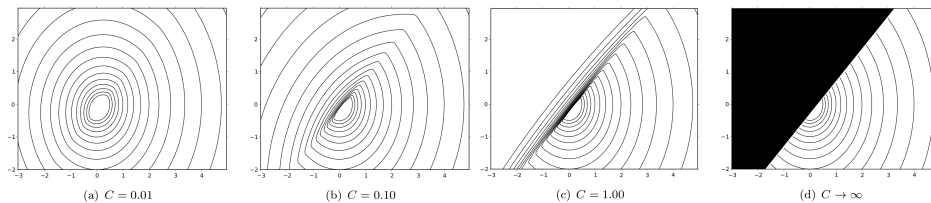
S-SVM: Toy example *

Consider a simple CRF model with a single variable, where $\mathcal{Y} = \{-1, +1\}$. We define the energy function as

$$E(x, y, \mathbf{w}) = w_1 \varphi_1(x, y) + w_2 \varphi_2(x, y) .$$

Assuming a training set $\mathcal{D} = \{(-10, +1), (-4, +1), (6, -1), (5, -1)\}$ with

$$\varphi_1(x, y) = \begin{cases} 0, & \text{if } y = -1 \\ x, & \text{if } y = +1 \end{cases} \quad \text{and} \quad \varphi_2(x, y) = \begin{cases} x, & \text{if } y = -1 \\ 0, & \text{if } y = +1 \end{cases} .$$

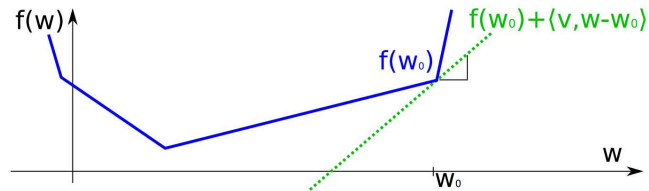


$$\frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{N} \sum_{n=1}^N \max_{\mathbf{y} \in \mathcal{Y}} (\Delta(\mathbf{y}^n, \mathbf{y}) - \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}) \rangle + \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}^n) \rangle) .$$

Subgradient *

Let $f : \mathbb{R}^D \rightarrow \mathbb{R}$ be a convex, but not necessarily differentiable, function. A vector $\mathbf{v} \in \mathbb{R}^D$ is called a **subgradient** of f at \mathbf{w}_0 , if

$$f(\mathbf{w}) \geq f(\mathbf{w}_0) + \langle \mathbf{v}, \mathbf{w} - \mathbf{w}_0 \rangle \quad \text{for all } \mathbf{w} .$$



Note that for differentiable f , the gradient $\mathbf{v} = \nabla f(\mathbf{w}_0)$ is the **only** subgradient.

Subgradient descent minimization *

Input: Tolerance $\epsilon > 0$ and step-sizes η_t .

Output: The minimizer \mathbf{w} of L .

- 1: $\mathbf{w} \leftarrow \mathbf{0}$
- 2: **repeat**
- 3: $\mathbf{v} \in \nabla_{\mathbf{w}}^{\text{sub}} L(\mathbf{w})$
- 4: $\mathbf{w} \leftarrow \mathbf{w} - \eta_t \mathbf{v}$
- 5: **until** L changed less than ϵ
- 6: **return** \mathbf{w}

This method converges to global minimum, but rather inefficient if the objective function L is non-differentiable.

Remark: For step sizes satisfying **diminishing step size conditions**:

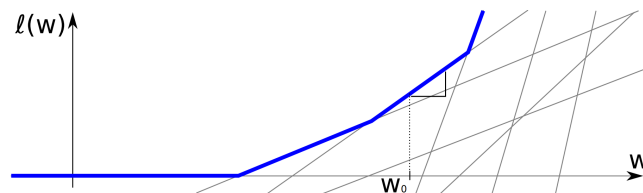
$$\lim_{t \rightarrow \infty} \eta_t = 0, \text{ and } \sum_{t=0}^{\infty} \eta_t \rightarrow \infty$$

convergence is guaranteed. For example $\eta_t := \frac{1+m}{t+m}$ for any $m \geq 0$.

Numerical solution

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{N} \sum_{n=1}^N \max_{\mathbf{y} \in \mathcal{Y}} (\Delta(\mathbf{y}^n, \mathbf{y}) - \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}) \rangle + \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}^n) \rangle).$$

As we have discussed, this function is non-differentiable. Therefore, we cannot use gradient descent directly, so we have to use subgradients.



For each $\mathbf{y} \in \mathcal{Y}$, ℓ is a linear function, since it is the maximum over all $\mathbf{y} \in \mathcal{Y}$. In order to calculate the subgradient at \mathbf{w}_0 , one may find the maximal (active) \mathbf{y} , and then use $\mathbf{v} = \nabla \ell(\mathbf{w}_0)$.

Calculating the subgradient

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{N} \sum_{n=1}^N \max_{\mathbf{y} \in \mathcal{Y}} (\Delta(\mathbf{y}^n, \mathbf{y}) - \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}) \rangle + \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}^n) \rangle).$$

Let $\bar{\mathbf{y}} \in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} \Delta(\mathbf{y}^n, \mathbf{y}) - \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}) \rangle$.

A subgradient \mathbf{v} is given by

$$\begin{aligned} & \nabla_{\mathbf{w}}^{\text{sub}} \left(\frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{N} \sum_{n=1}^N \max_{\mathbf{y} \in \mathcal{Y}} (\Delta(\mathbf{y}^n, \mathbf{y}) - \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}) \rangle + \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}^n) \rangle) \right) \\ & \ni \nabla_{\mathbf{w}} \left(\frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{N} \sum_{n=1}^N (\Delta(\mathbf{y}^n, \bar{\mathbf{y}}) - \langle \mathbf{w}, \varphi(\mathbf{x}^n, \bar{\mathbf{y}}) \rangle + \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}^n) \rangle) \right) \\ & = \mathbf{w} + \frac{C}{N} \sum_{n=1}^N -\varphi(\mathbf{x}^n, \bar{\mathbf{y}}) + \varphi(\mathbf{x}^n, \mathbf{y}^n) =: \mathbf{v}. \end{aligned}$$

Subgradient descent S-SVM learning *

Input: Training set $\mathcal{D} = \{(\mathbf{x}^1, \mathbf{y}^1), \dots, (\mathbf{x}^n, \mathbf{y}^n)\}$, energies $\varphi(\mathbf{x}, \mathbf{y})$, loss function $\Delta(\mathbf{y}, \mathbf{y}')$, regularizer C and step-sizes η_1, \dots, η_T for all the T iterations.

Output: the weight vector \mathbf{w} for the prediction function $f(\mathbf{x}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} -\langle \mathbf{w}, \varphi(\mathbf{x}, \mathbf{y}) \rangle$.

```
1:  $\mathbf{w} \leftarrow \mathbf{0}$ 
2: for  $t = 1, \dots, T$  do
3:   for  $n = 1, \dots, N$  do
4:      $\bar{\mathbf{y}} \leftarrow \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} \Delta(\mathbf{y}^n, \mathbf{y}) - \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}) \rangle$ 
5:      $\mathbf{v}^n \leftarrow -\varphi(\mathbf{x}^n, \bar{\mathbf{y}}) + \varphi(\mathbf{x}^n, \mathbf{y}^n)$ 
6:   end for
7:    $\mathbf{w} \leftarrow \mathbf{w} - \eta_t \left( \underbrace{\mathbf{w} + \frac{C}{N} \sum_{n=1}^N \mathbf{v}^n}_{\mathbf{v}} \right)$ 
8: end for
```

The step-size can be chosen as $\eta_t = \frac{1}{t}$ for all $t = 1, \dots, T$.

Note that each update of \mathbf{w} needs only one argmax-prediction.

Stochastic subgradient descent S-SVM learning *

Input: Training set $\mathcal{D} = \{(\mathbf{x}^1, \mathbf{y}^1), \dots, (\mathbf{x}^n, \mathbf{y}^n)\}$, energies $\varphi(\mathbf{x}, \mathbf{y})$, loss function $\Delta(\mathbf{y}, \mathbf{y}')$, regularizer C and step-sizes η_1, \dots, η_T for all the T iterations.

Output: The weight vector w for the prediction function $f(\mathbf{x}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} -\langle \mathbf{w}, \varphi(\mathbf{x}, \mathbf{y}) \rangle$.

```
1:  $\mathbf{w} \leftarrow \mathbf{0}$ 
2: for  $t = 1, \dots, T$  do
3:    $(\mathbf{x}^n, \mathbf{y}^n) \leftarrow$  a randomly chosen training example
4:    $\bar{\mathbf{y}} \leftarrow \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} \Delta(\mathbf{y}^n, \mathbf{y}) - \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}) \rangle$ 
5:    $\mathbf{w} \leftarrow \mathbf{w} - \eta_t \left( \mathbf{w} + \frac{C}{N} (-\varphi(\mathbf{x}^n, \bar{\mathbf{y}}) + \varphi(\mathbf{x}^n, \mathbf{y}^n)) \right)$ 
6: end for
```

Note that each update step of w needs only one argmax-prediction, however we will generally need **many** iterations until convergence.

Summary of S-SVM learning *

We are given a *training set* $\mathcal{D} = \{(\mathbf{x}^1, \mathbf{y}^1), \dots, (\mathbf{x}^n, \mathbf{y}^n)\} \subset \mathcal{X} \times \mathcal{Y}$ and a problem specific *loss function* $\Delta : \mathcal{Y} \times \mathcal{Y} \leftarrow \mathbb{R}$.

The task is to *learn* parameter \mathbf{w} for *prediction function*

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

that minimizes *expected loss* on the *training set*.

S-SVM solution derived by *maximum margin framework*:

$$\langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}) \rangle \leq \langle \mathbf{w}, \varphi(\mathbf{x}^n, \mathbf{y}^n) \rangle + \Delta(\mathbf{y}^n, \mathbf{y}) ,$$

that is the predicted output is enforced to be not worse than the correct one by a *margin*.

We have seen that S-SVM learning ends up a **convex optimization** problem, but it is **non-differentiable**. Furthermore it requires repeated *argmax prediction*.

Next lecture: Summary of the course *



Literature *

1. Sebastian Nowozin and Christoph H. Lampert. Structured prediction and learning in computer vision. *Foundations and Trends in Computer Graphics and Vision*, 6(3–4), 2010
2. Christopher Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006