# Loss Functions for Deep Structured Models

Presenter: Jack Lanchantin

University of Virginia
https://qdata.github.io/deep2Read/

12/12/18

# Outline

# SVM Losses

# SVM Loss for Binary Classification

$D = \{(x_i, y_i)\}_{i=1}^{n}, x_i \in \mathbb{R}^d, y_i \in \{0, 1\}$

$$\mathcal{L}(w) = \frac{1}{n} \sum_{i=1}^{n} \max(0, 1 - y_i \langle w, x_i \rangle)$$

If $y_i = 1$, $\langle w, x_i \rangle$ should be $\geq 1$
If $y_i = -1$, $\langle w, x_i \rangle$ should be $\leq -1$
both of these result in $1 - 1 = 0$, giving an SVM loss of 0.

# SVM Loss for Multi-Class Classification

$D = \{(x_i, y_i)\}_{i=1}^{n}, \; x_i \in \mathbb{R}^d, \; y_i \in \{0, 1, 2, ..., c\}$

$$\mathcal{L}(w) = \frac{1}{n} \sum_{i=1}^{n} \max(0, 1 + \max_{t \neq y_i} \langle w_t, x_i \rangle - \langle w_{y_i}, x_i \rangle)$$

This maximizes the margin between the true label $y_i$ and the next biggest label's prediction ($\max_{t \neq y_i}$).

Alternatively, we could push down *all* labels which aren't the true label (as opposed to only the next max):

$$\mathcal{L}(w) = \frac{1}{n} \sum_{i=1}^{n} \sum_{t \neq y_i} \max(0, 1 + \langle w_t, x_i \rangle - \langle w_{y_i}, x_i \rangle)$$

# Structured SVM for Multi-label Classification

$D = \{(x_i, y_i)\}_{i=1}^n$, $x_i \in \mathbb{R}^d$, $y_i \in \{0, 1\}^L$, where $L$ the number of labels. $\hat{y}$ is the predicted output set, and $\Delta$ is hamming loss

$$\mathcal{L}(w) = \frac{1}{n} \sum_{i=1}^n \max(0, \max_{\hat{y}} (\Delta(y, \hat{y}) + \langle w, \phi(x_i, \hat{y}) \rangle) - \langle w, \phi(x_i, y) \rangle)$$

Finding this $\max_{\hat{y}}$ is typically intractable, but gradient descend can be used to approximate.

# SPEN and SPEN InfNet

## Energy Models

Instead of maximizing the score of $\langle w, \phi(x_i, \hat{y}) \rangle$, we can instead minimize the energy of the joint pair $(x, y)$

$$\min_y E(x, y) \quad \text{subject to} \quad y \in \{0, 1\}^L$$

This could be rendered tractable by assuming certain structure (e.g., a tree) for the energy function $E(x, y)$. Instead, we consider a general $E(x, y)$, but optimize over a convex relaxation of the constraint set

$$\min_y E(x, y) \quad \text{subject to} \quad y \in [0, 1]^L$$

# Structured Prediction Energy Networks (SPEN)

The energy of a pair $(x, y)$ is represented as:

$$E(x, y) = E(x, y)^{local} + E(x, y)^{label}$$

$$E(x, y)^{local} = \sum_{i=1}^{L} y_i b_i^\top F(x)$$

$$E(x, y)^{label} = c_2^\top \sigma(C_1 y)$$

# Structured Prediction Energy Networks (SPEN)

If we were to use the CRF framework in SPENs, the label, or global energy network would be

$$\mathrm{E}(\mathrm{x}, \mathrm{y})^{label} = \mathrm{y}^{\top} S_1 \mathrm{y} + s \top \mathrm{y}$$

which doesn't use the $[0, 1]^L$ relaxation on y and only considers pairwise dependencies as opposed to the $C_1$ matrix in the original label network, which allows for any learned dependencies.

# Structured Prediction Energy Networks (SPEN)

SPENs minimize the following SSVM Loss, where $[\cdot]_+$ represents $\max(0, \cdot)$

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^{n} \left[ \max_{\hat{y}} (\Delta(y_i, \hat{y}) - E(x_i, \hat{y}) + E(x_i, y_i)) \right]_+$$

During training, the $\max_{\hat{y}}$ is found by *cost-augmented inference*, where $\mathcal{P}$ projects values into [0,1]:
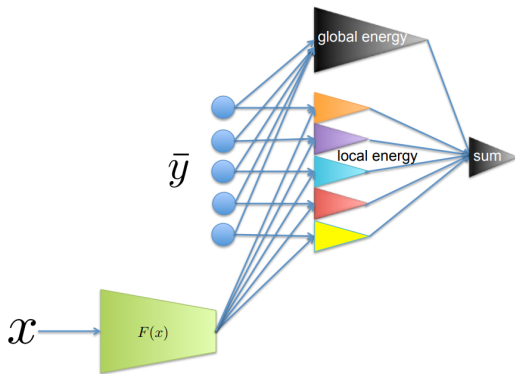
$$y^{t+1} = \mathcal{P}\left(y^t - \eta \frac{d}{dy}(-\Delta(y_i, y^t) + E(x, y^t))\right) \qquad (1)$$

# Structured Prediction Energy Networks (SPEN)

At test time, given $x_i$, $\hat{y}$ is found by minimizing the energy $\mathrm{E}(x_i, \hat{y})$. This is done via gradient descent:

$$y^{t+1} = \mathcal{P}\Big(y^t - \eta \frac{\mathrm{d}}{\mathrm{dy}}\mathrm{E}(x, y^t))\Big)$$

# SPEN Model

# Learning Approximate Inference Networks for Structured Prediction (SPEN InfNet)

Instead of using (13), learn an *inference network* $G_\Psi$ with the goal that

$$G_\Psi(x) \approx \arg\min_{\hat{y}} E(x, \hat{y})$$

Given energy function $E$, we seek to minimize the following:

$$\arg\min_\Psi E(x, G_\Psi(x))$$

# Learning Approximate Inference Networks for Structured Prediction (SPEN InfNet)

Training InfNet requires two steps: InfNet turns the loss into a minimax problem to learn a cost augmented inference network $G_\Phi$

$$\min_\Theta \max_\Phi \frac{1}{n} \sum_{i=1}^n \Big[ \Delta(y_i, G_\Phi(x_i)) - \mathrm{E}_\Theta(x_i, G_\Phi(x_i)) + \mathrm{E}_\Theta(x_i, y_i) \Big]_+$$

$$\hat{\Theta} \leftarrow \arg\min_\Theta \Big[ \Delta(y_i, G_\Phi(x_i)) - \mathrm{E}_\Theta(x_i, G_\Phi(x_i)) + \mathrm{E}_\Theta(x_i, y_i) \Big]_+$$

$$\hat{\Phi} \leftarrow \arg\min_\Phi \Big[ -\Delta(y_i, G_\Phi(x_i)) + \mathrm{E}_\Theta(x_i, G_\Phi(x_i)) - \mathrm{E}_\Theta(x_i, y_i) \Big]_+$$

# Learning Approximate Inference Networks for Structured Prediction (SPEN InfNet)

Training only gives us a cost-augmented inference network $G_\Phi$, but we want inference network $G_\Psi$. So we initialize $G_\Psi$ with $G_\Phi$ and minimize the original Energy:

$$\arg\min_\Psi E(x, G_\Psi(x))$$

# Deep Value Networks

# Deep Value Networks

Train a value function $v(x, \hat{y}; \theta)$ that approximates an oracle function $v^*$ (dependent on task):

$$v(x, \hat{y}; \theta) \approx v^*(y, \hat{y})$$

For MLC,

$$v^*(y, \hat{y}) = v_{F1}(y, \hat{y}) = \frac{2(y \cap \hat{y})}{(y \cap \hat{y}) + (y \cup \hat{y})}$$

At inference time, $\hat{y}$ is found by initializing to 0 vector $\hat{y}^0 = [0]^L$ and then updating via gradient ascent on $v(x, \hat{y}^t; \theta)$

$$\hat{y}^{t+1} = \mathcal{P}\left(\hat{y}^t + \eta \frac{\mathrm{d}}{\mathrm{d}\hat{y}} v(x, \hat{y}^t; \theta)\right)$$

# Deep Value Networks

Our training objective aims at minimizing the discrepancy between $v(x, \hat{y}; \theta)$ and $v^*(y, \hat{y})$ on a training set of triplets (input, output, value) denoted $D = \{x^i, y^i, v^{*i}\}_{i=1}^N$

This can be done using binary cross entropy between $v(x, \hat{y}; \theta)$ and $v^*(y, \hat{y})$:

$$\mathcal{L} = \sum_{i=1}^n -v^*(y_i, \hat{y}) \log(v(x_i, \hat{y}; \theta)) - (1 - v^*(y_i, \hat{y})) \log(1 - v(x_i, \hat{y}; \theta))$$

# Deep Value Networks

Choosing $\hat{y}$ can be done in several ways:

- **Ground Truth**: $\hat{y} = y^i$ and $v * (\hat{y}, y^i) = 1$
- **Inference**: Gradient ascent on the current $v^*$
- **Adversarial Samples**: Maximize the cross entropy loss using gradient ascent

Other

# Adversarial Training for Segmentation

$$\hat{\Theta} \leftarrow \underset{\Theta}{\arg\min} \left( \ell_{BCE}(D(\mathsf{x}_i, G_\Phi(\mathsf{x}_i), 0) + \ell_{BCE} D(\mathsf{x}_i, \mathsf{y}_i), 1) \right)$$

$$\hat{\Phi} \leftarrow \underset{\Phi}{\arg\min} \left( -\ell_{BCE}(D(\mathsf{x}_i, G_\Phi(\mathsf{x}_i), 0) + \ell_{avgBCE}(\mathsf{y}, G_\Phi(\mathsf{x}_i)) \right)$$

# Other Loss Functions

$N$ samples, $L$ labels. True label $y_j^i$ and predicted label $\hat{y}_j^i$ for sample $i$ and label $j$.

## BCE

$$\mathcal{L}_{BCE} = \sum_{i=1}^{N} \sum_{j=1}^{L} -\left(y_j^i \log \hat{y}_j^i + (1 - y_j^i) \log(1 - \hat{y}_j^i)\right) \qquad (2)$$

## L2

$$\mathcal{L}_{L2} = \sum_{i=1}^{N} \sum_{j=1}^{L} (y_j^i - \hat{y}_j^i)^2 \qquad (3)$$

## KL

$$\mathcal{L}_{KL} = \sum_{i=1}^{N} \sum_{j=1}^{L} y_j^i \log \left(\frac{y_j^i}{\hat{y}_j^i}\right) \qquad (4)$$