Empirical Risk Minimization

Outline

- Empirical risk minimization view
 - Perceptron
 - CRF

Warning: Math Ahead

Notation for Linear Models

- Training data: {(x_1, y_1), (x_2, y_2), ..., (x_N, y_N)}
- Testing data: {(x_{N+1}, y_{N+1}), ... ($x_{N+N'}, y_{N+N'}$)}
- Feature function: **g**
- Weights: w
- Decoding:

$$ext{decode}(\mathbf{w}, oldsymbol{x}) = rg\max_{oldsymbol{y}} \mathbf{w}^{ op} \mathbf{g}(oldsymbol{x}, oldsymbol{y})$$

- Learning: $\operatorname{learn}\left(\{(\boldsymbol{x}_i, \boldsymbol{y}_i)\}_{i=1}^N\right) = \operatorname{arg}\max_{\boldsymbol{w}} \Phi\left(\boldsymbol{w}, \{(\boldsymbol{x}_i, \boldsymbol{y}_i)\}_{i=1}^N\right)$
- Evaluation:

$$\frac{1}{N'} \sum_{i=1}^{N'} \operatorname{cost} \left(\operatorname{decode} \left(\operatorname{learn} \left(\{ (\boldsymbol{x}_i, \boldsymbol{y}_i) \}_{i=1}^N \right), \boldsymbol{x}_{N+i} \right), \boldsymbol{y}_{N+i} \right)$$

Structured Perceptron

- Described as an online algorithm.
- On each iteration, take one example, and update the weights according to:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \left(\mathbf{g}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{g}(\mathbf{x}_t, \operatorname{decode}(\mathbf{w}, \mathbf{x}_t)) \right)$$

 Not discussing today: the theoretical guarantees this gives, separability, and the averaged and voted versions.

Empirical Risk Minimization

• A unifying framework for many learning algorithms.

$$\operatorname{learn}\left(\{(\boldsymbol{x}_{i}, \boldsymbol{y}_{i})\}_{i=1}^{N}\right) = \operatorname{arg\,max}_{\boldsymbol{w}} \Phi\left(\boldsymbol{w}, \{(\boldsymbol{x}_{i}, \boldsymbol{y}_{i})\}_{i=1}^{N}\right)$$
$$= \operatorname{arg\,min}_{\boldsymbol{w}} \underbrace{\frac{1}{N} \sum_{i=1}^{N} L(\boldsymbol{w}, \boldsymbol{x}_{i}, \boldsymbol{y}_{i})}_{\approx \mathbb{E}[L(\boldsymbol{w}, \boldsymbol{X}, \boldsymbol{Y})]} + R(\boldsymbol{w})$$

• Many options for the loss function L and the regularization function R.

Loss Functions You May Know

Name	Expression of $L(\mathbf{w}, oldsymbol{x}, oldsymbol{y})$
Log loss (joint)	$-\log p(\boldsymbol{x}, \boldsymbol{y} \mid \mathbf{w})$
Log loss (conditional)	$-\log p(\boldsymbol{y} \mid \boldsymbol{x}, \mathbf{w})$
Zero-one loss	$1\{\operatorname{decode}(\mathbf{w}, \boldsymbol{x}) \neq \boldsymbol{y}\}$
Expected zero-one loss	$1 - p(\boldsymbol{y} \mid \boldsymbol{x}, \mathbf{w})$

Loss Functions You May Know

Name	Expression of $L(\mathbf{w}, oldsymbol{x}, oldsymbol{y})$
Log loss (joint)	$-\log p(\boldsymbol{x}, \boldsymbol{y} \mid \mathbf{w})$
Log loss (conditional)	$-\log p(\boldsymbol{y} \mid \boldsymbol{x}, \mathbf{w})$
Cost	$\operatorname{cost}(\operatorname{decode}(\mathbf{w}, \boldsymbol{x}), \boldsymbol{y})$
Expected cost, a.k.a. "risk"	$\mathbb{E}_{p(oldsymbol{Y} oldsymbol{x},oldsymbol{w})}[ext{cost}(oldsymbol{Y},oldsymbol{y})]$

Some Questions

- Where do CRFs fit into this picture?
- Where does the structured perceptron fit into this picture?
- Do people directly minimize cost or expected cost?

CRFs and Loss

 Plugging in the log-linear form (and not worrying at this level about locality of features):

$$\begin{aligned} -\log p(\boldsymbol{y} \mid \boldsymbol{x}, \mathbf{w}) &= -\mathbf{w}^{\top} \mathbf{g}(\boldsymbol{x}, \boldsymbol{y}) + \log \sum_{\boldsymbol{y}'} \exp \mathbf{w}^{\top} \mathbf{g}(\boldsymbol{x}, \boldsymbol{y}) \\ \frac{\partial L}{\partial w_{j}} &= -g_{j}(\boldsymbol{x}, \boldsymbol{y}) + \mathbb{E}_{p(\boldsymbol{Y} \mid \boldsymbol{x}, \mathbf{w})}[g_{j}(\boldsymbol{x}, \boldsymbol{Y})] \end{aligned}$$

Training CRFs and Other Log-Linear Models

- Early days: iterative scaling (specialized method)
- ~2002: quasi-Newton methods
- ~2006: stochastic gradient descent (LeCun, 1998)

Perceptron and Loss

 Not clear immediately what L is, but the "gradient" of L should be:

$$-g_j(\boldsymbol{x}, \boldsymbol{y}) + g_j(\boldsymbol{x}, \operatorname{decode}(\mathbf{w}, \boldsymbol{x}))$$

• The vector of above quantities is actually a *subgradient* of:

$$L(\mathbf{w}, \boldsymbol{x}, \boldsymbol{y}) = -\mathbf{w}^{\top} \mathbf{g}(\boldsymbol{x}, \boldsymbol{y}) + \max_{\boldsymbol{y}'} \mathbf{w}^{\top} \mathbf{g}(\boldsymbol{x}, \boldsymbol{y}')$$

Compare

- CRF (log-loss):
- $-\log p(\boldsymbol{y} \mid \boldsymbol{x}, \mathbf{w}) = -\mathbf{w}^{\top} \mathbf{g}(\boldsymbol{x}, \boldsymbol{y}) + \log \sum_{\boldsymbol{y}'} \exp \mathbf{w}^{\top} \mathbf{g}(\boldsymbol{x}, \boldsymbol{y})$
 - Perceptron:

$$L(\mathbf{w}, \boldsymbol{x}, \boldsymbol{y}) = -\mathbf{w}^{\top} \mathbf{g}(\boldsymbol{x}, \boldsymbol{y}) + \max_{\boldsymbol{y}'} \mathbf{w}^{\top} \mathbf{g}(\boldsymbol{x}, \boldsymbol{y}')$$

Loss Functions You Know

Name	Expression of $L(\mathbf{w}, oldsymbol{x}, oldsymbol{y})$	Convex?
Log loss (joint)	$-\log p(oldsymbol{x},oldsymbol{y}\mid\mathbf{w})$	✓
Log loss (conditional)	$-\log p(oldsymbol{y} \mid oldsymbol{x}, \mathbf{w})$	~
Cost	$\mathrm{cost}(\mathrm{decode}(\mathbf{w}, \boldsymbol{x}), \boldsymbol{y})$	
Expected cost, a.k.a. "risk"	$\mathbb{E}_{p(oldsymbol{Y} oldsymbol{x},oldsymbol{w})}[ext{cost}(oldsymbol{Y},oldsymbol{y})]$	
Perceptron loss	$-\mathbf{w}^{ op}\mathbf{g}(\boldsymbol{x}, \boldsymbol{y}) + \max_{\boldsymbol{y}'} \mathbf{w}^{ op}\mathbf{g}(\boldsymbol{x}, \boldsymbol{y}')$	~

Loss Functions You Know

Name	Expression of $L(\mathbf{w}, oldsymbol{x}, oldsymbol{y})$	Cont.?
Log loss (joint)	$-\log p(oldsymbol{x},oldsymbol{y}\mid\mathbf{w})$	✓
Log loss (conditional)	$-\log p(oldsymbol{y} \mid oldsymbol{x}, \mathbf{w})$	~
Cost	$\mathrm{cost}(\mathrm{decode}(\mathbf{w}, \boldsymbol{x}), \boldsymbol{y})$	
Expected cost, a.k.a. "risk"	$\mathbb{E}_{p(oldsymbol{Y} oldsymbol{x}, oldsymbol{w})}[ext{cost}(oldsymbol{Y}, oldsymbol{y})]$	✓
Perceptron loss	$-\mathbf{w}^{ op}\mathbf{g}(\boldsymbol{x}, \boldsymbol{y}) + \max_{\boldsymbol{y}'} \mathbf{w}^{ op}\mathbf{g}(\boldsymbol{x}, \boldsymbol{y}')$	~

Loss Functions You Know

Name	Expression of $L(\mathbf{w}, oldsymbol{x}, oldsymbol{y})$	Cost?
Log loss (joint)	$-\log p(\boldsymbol{x}, \boldsymbol{y} \mid \mathbf{w})$	
Log loss (conditional)	$-\log p(oldsymbol{y} \mid oldsymbol{x}, \mathbf{w})$	
Cost	$\mathrm{cost}(\mathrm{decode}(\mathbf{w}, \boldsymbol{x}), \boldsymbol{y})$	~
Expected cost, a.k.a. "risk"	$\mathbb{E}_{p(oldsymbol{Y} oldsymbol{x},oldsymbol{w})}[ext{cost}(oldsymbol{Y},oldsymbol{y})]$	~
Perceptron loss	$-\mathbf{w}^{ op}\mathbf{g}(oldsymbol{x},oldsymbol{y}) + \max_{oldsymbol{y}'}\mathbf{w}^{ op}\mathbf{g}(oldsymbol{x},oldsymbol{y}')$	

The Ideal Loss Function

For computational convenience:

- Convex
- Continuous
- For good performance:
- Cost-aware
- Theoretically sound

On Regularization

- In principle, this choice is orthogonal to the loss function.
- Squared L₂ norm is the most common starting place.

$$R(\mathbf{w}) = \mathbf{k} \|\mathbf{w}\|_2^2$$

$$= \lambda \sum_{j} w_{j}^{2}$$

• L₁ and other sparsity-inducing $R(\mathbf{w}) = \lambda \|\mathbf{w}\|_1$ regularizers are attracting more attention lately. $= \lambda \sum |i|$

$$= \lambda \sum_{j} |w_{j}|$$

Practical Advice

- Features still more important than the loss function.
 - But general, easy-to-implement algorithms are quite useful!
- Perceptron is easiest to implement.
- CRFs and max margin techniques usually do better.
- Tune the regularization constant, λ .
 - Never on the test data.