

# Empirical Risk Minimization

October 29, 2015

# Outline

- Empirical risk minimization view
  - Perceptron
  - CRF

# Notation for Linear Models

- Training data:  $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$
- Testing data:  $\{(\mathbf{x}_{N+1}, y_{N+1}), \dots, (\mathbf{x}_{N+N'}, y_{N+N'})\}$
- Feature function:  $\mathbf{g}$
- Weights:  $\mathbf{w}$
- Decoding:
$$\text{decode}(\mathbf{w}, \mathbf{x}) = \arg \max_{\mathbf{y}} \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y})$$
- Learning:
$$\text{learn}(\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N) = \arg \max_{\mathbf{w}} \Phi(\mathbf{w}, \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N)$$
- Evaluation:
$$\frac{1}{N'} \sum_{i=1}^{N'} \text{cost}(\text{decode}(\text{learn}(\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N), \mathbf{x}_{N+i}), \mathbf{y}_{N+i})$$

# Structured Perceptron

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

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

- 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.

$$\begin{aligned} \text{learn} \left( \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \right) &= \arg \max_{\mathbf{w}} \Phi \left( \mathbf{w}, \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \right) \\ &= \arg \min_{\mathbf{w}} \underbrace{\frac{1}{N} \sum_{i=1}^N L(\mathbf{w}, \mathbf{x}_i, \mathbf{y}_i) + R(\mathbf{w})}_{\approx \mathbb{E}[L(\mathbf{w}, \mathbf{X}, \mathbf{Y})]} \end{aligned}$$

- Many options for the loss function  $L$  and the regularization function  $R$ .

# Solving the Minimization Problem

- In some friendly cases, there is a closed form solution for the minimizer of  $w$ 
  - E.g., the maximum likelihood estimator for HMMs
- *Usually*, we have to use an iterative algorithm which amounts to progressively finding better versions of  $w$ 
  - involves hard/soft inference with each improved value of  $w$  on either part or all of the training set

# Loss Functions You May Know

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

# Loss Functions You May Know

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



# Loss Functions You May Know

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

# CRFs and Loss

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

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

$$\frac{\partial L}{\partial w_j} = -g_j(\mathbf{x}, \mathbf{y}) + \mathbb{E}_{p(\mathbf{Y} \mid \mathbf{x}, \mathbf{w})} [g_j(\mathbf{x}, \mathbf{Y})]$$

# CRFs and Loss

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

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

$$\frac{\partial L}{\partial w_j} = -g_j(\mathbf{x}, \mathbf{y}) + \mathbb{E}_{p(\mathbf{Y} \mid \mathbf{x}, \mathbf{w})} [g_j(\mathbf{x}, \mathbf{Y})]$$

$$\frac{\partial L}{\partial \mathbf{w}} = -\mathbf{g}(\mathbf{x}, \mathbf{y}) + \mathbb{E}_{p(\mathbf{Y} \mid \mathbf{x}, \mathbf{w})} [\mathbf{g}(\mathbf{x}, \mathbf{Y})]$$

# Training CRFs and Other Linear Models

- Early days: iterative scaling (specialized method for log-linear models only)
- ~2002: quasi-Newton methods
  - (using LBFGS which dates from the late 1980s)
- ~2006: stochastic gradient descent
- ~2010: adaptive gradient methods

# Perceptron and Loss

- Not clear immediately what  $L$  is, but the “gradient” of  $L$  should be:

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

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

$$L(\mathbf{w}, \mathbf{x}, \mathbf{y}) = \max\{0, -\mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y})\} + \max_{\hat{\mathbf{y}}} \mathbf{w}^\top (\mathbf{x}, \hat{\mathbf{y}})$$

# Compare

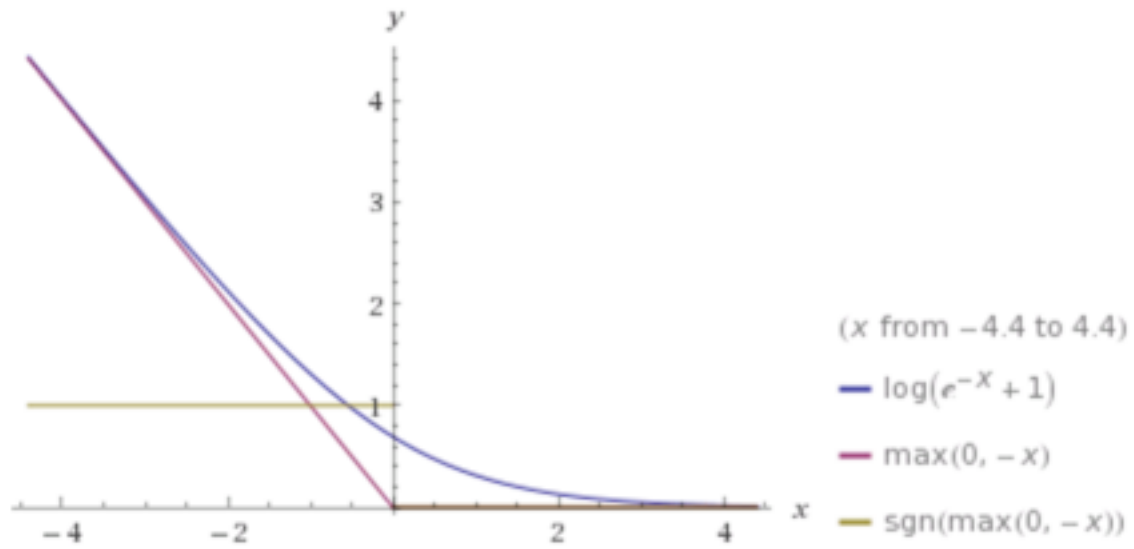
- CRF (log-loss):

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

- Perceptron:

$$L(\mathbf{w}, \mathbf{x}, \mathbf{y}) = \max\{0, -\mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}) + \max_{\hat{\mathbf{y}}} \mathbf{w}^\top (\mathbf{x}, \hat{\mathbf{y}})\}$$

# Loss Functions



# Loss Functions You Know

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



# Loss Functions You Know

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

# Loss Functions You Know

Name	Expression of $L(\mathbf{w}, \mathbf{x}, \mathbf{y})$	Cost?
Log loss (joint)	$-\log p(\mathbf{x}, \mathbf{y} \mid \mathbf{w})$	
Log loss (conditional)	$-\log p(\mathbf{y} \mid \mathbf{x}, \mathbf{w})$	
Cost	$\text{cost}(\text{decode}(\mathbf{w}, \mathbf{x}), \mathbf{y})$	✓
Expected cost, a.k.a. “risk”	$\mathbb{E}_{p(\mathbf{Y} \mid \mathbf{x}, \mathbf{w})} [\text{cost}(\mathbf{Y}, \mathbf{y})]$	✓
Perceptron loss	$-\mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}) + \max_{\mathbf{y}'} \mathbf{w}^\top \mathbf{g}(\mathbf{x}, \mathbf{y}')$	

# The Ideal Loss Function

For computational convenience:

- Convex
- Continuous

For good performance:

- Cost-aware
- Theoretically sound

# On Regularization

- In principle, this choice is independent from the choice of the loss function.
- Squared  $L_2$  norm is the most common starting place.

$$\begin{aligned} R(\mathbf{w}) &= \lambda \|\mathbf{w}\|_2^2 \\ &= \lambda \sum_j w_j^2 \end{aligned}$$

- $L_1$  and other sparsity-inducing regularizers as well as structured **regularizers** are of interest

$$\begin{aligned} R(\mathbf{w}) &= \lambda \|\mathbf{w}\|_1 \\ &= \lambda \sum_j |w_j| \end{aligned}$$

# 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,  $\lambda$ .
  - Never on the test data.