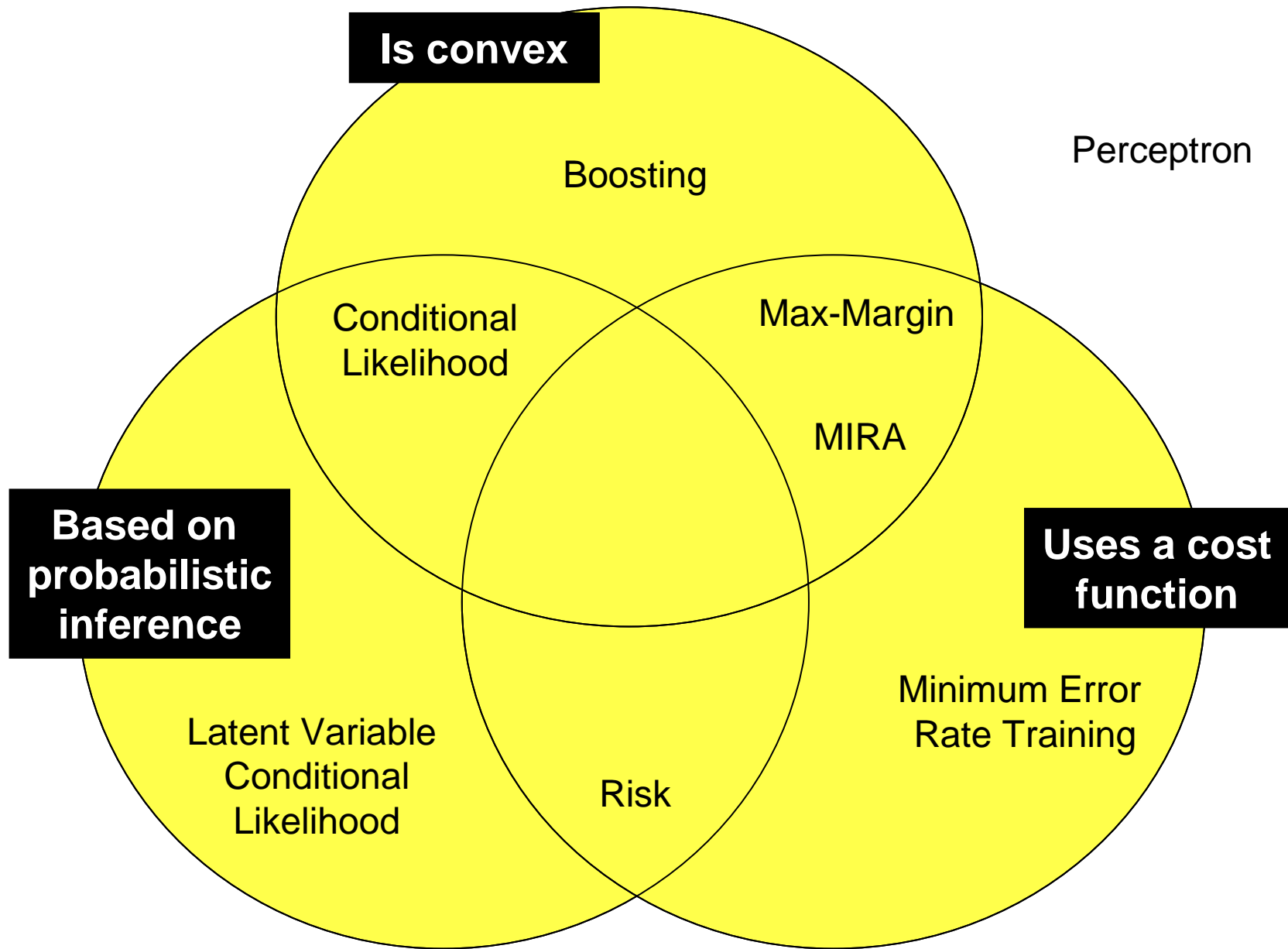


Softmax-Margin CRFs: Training Log-Linear Models with Cost Functions

Kevin Gimpel and Noah A. Smith



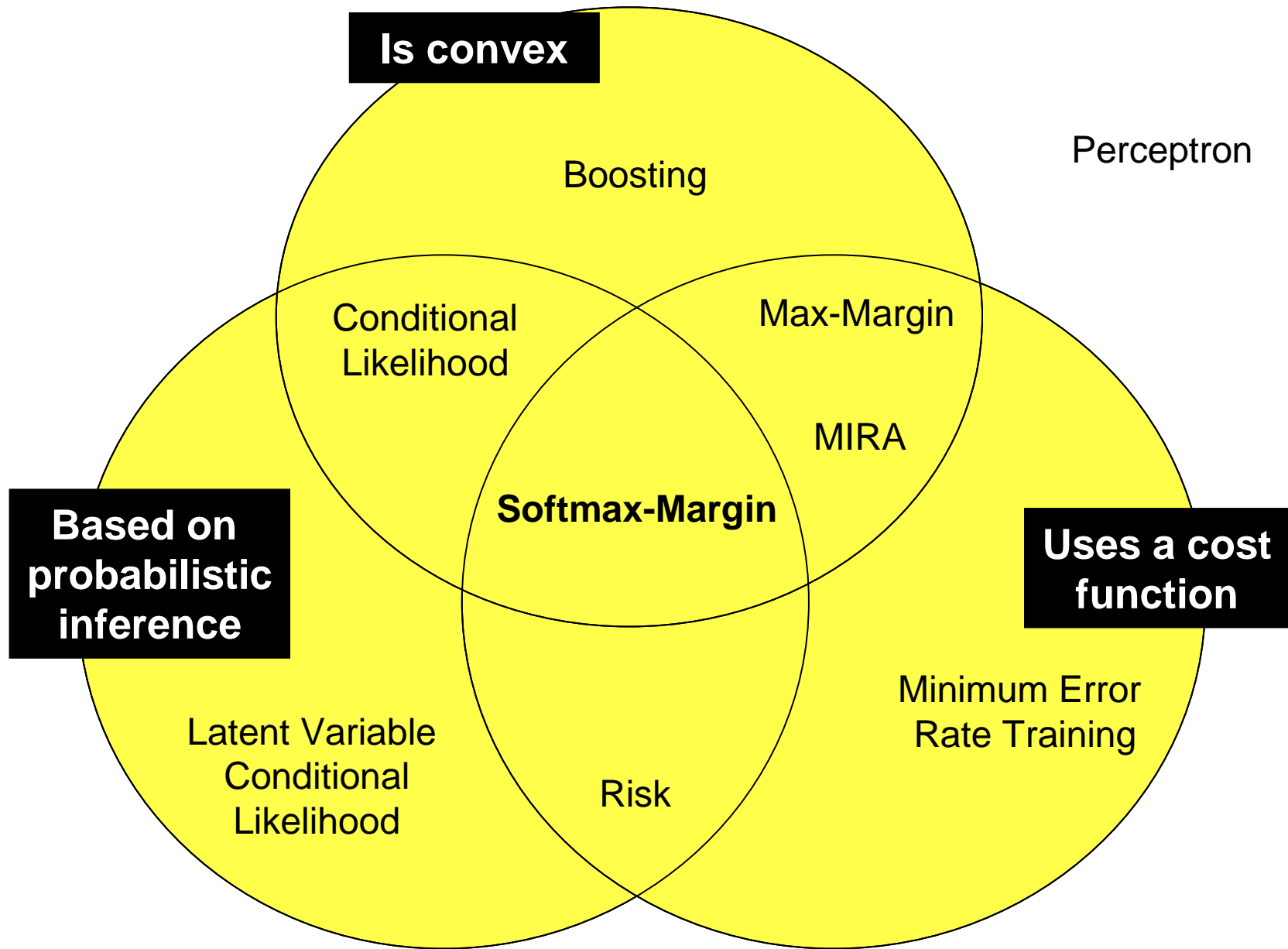
Carnegie Mellon



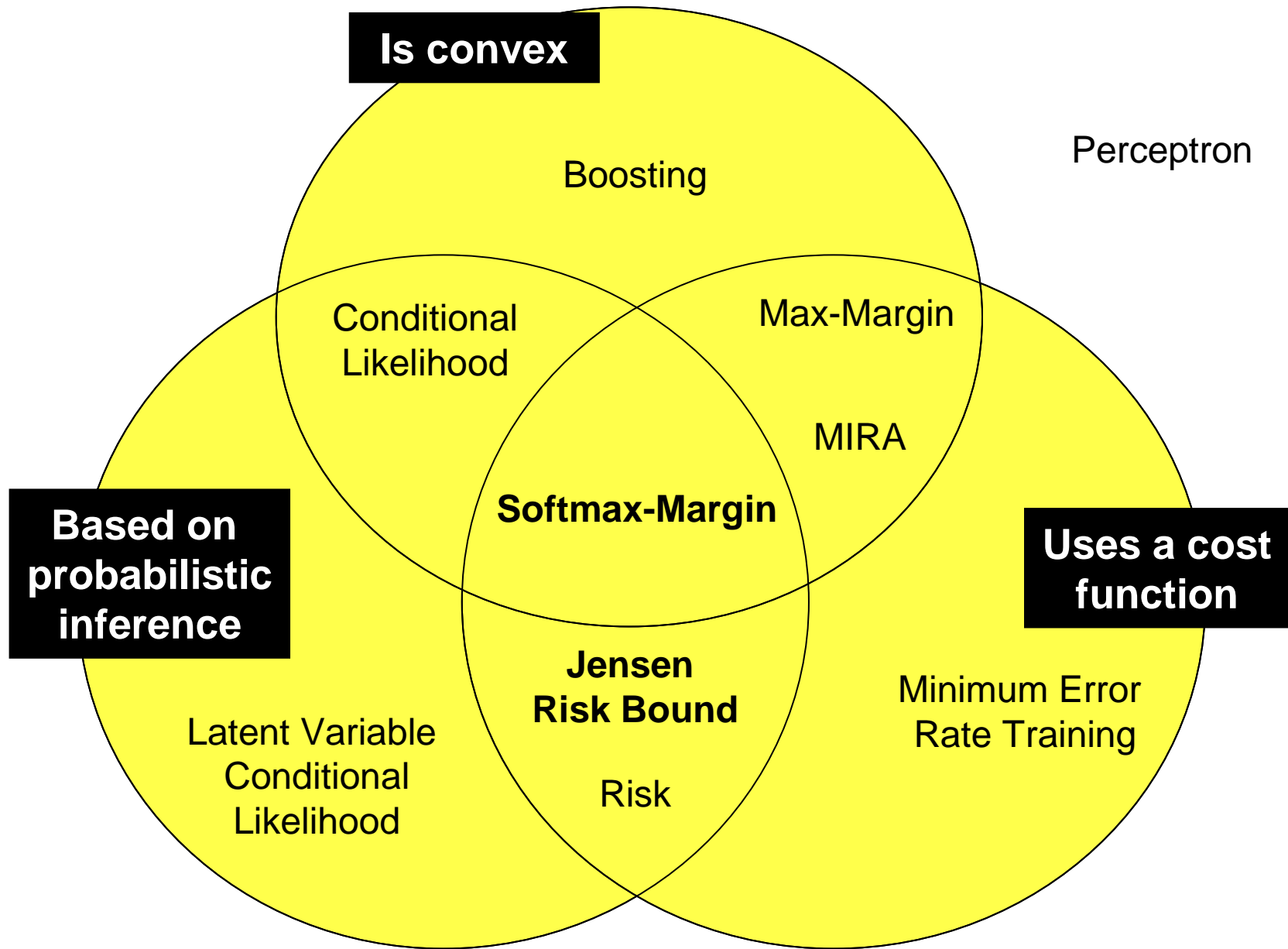
Perceptron



Carnegie Mellon



Carnegie Mellon



Carnegie Mellon

Linear Models for Structured Prediction

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}(x)} \boldsymbol{\theta}^\top \mathbf{f}(x, y)$$

input → x *output* → y
weights → $\boldsymbol{\theta}$ *features* → $\mathbf{f}(x, y)$

- For probabilistic interpretation, exponentiate and normalize:

$$p_{\boldsymbol{\theta}}(y|x) = \frac{\exp\{\boldsymbol{\theta}^\top \mathbf{f}(x, y)\}}{\sum_{y' \in \mathcal{Y}(x)} \exp\{\boldsymbol{\theta}^\top \mathbf{f}(x, y')\}}$$



Carnegie Mellon

Training

- Standard approach is to maximize conditional likelihood:

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^n \left(-\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y^{(i)}) + \log \sum_{y \in \mathcal{Y}(x^{(i)})} \exp\{\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y)\} \right)$$

- Another approach maximizes margin (Taskar et al., 2003):

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^n \left(-\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y^{(i)}) + \max_{y \in \mathcal{Y}(x^{(i)})} \left(\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y) + \text{cost}(y^{(i)}, y) \right) \right)$$

*task-specific
cost function*



Carnegie Mellon

Training

- Standard approach is to maximize conditional likelihood:

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^n \left(-\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y^{(i)}) + \log \sum_{y \in \mathcal{Y}(x^{(i)})} \exp\{\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y)\} \right)$$

- Another approach maximizes margin (Taskar et al., 2003):

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^n \left(-\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y^{(i)}) + \underbrace{\max_{y \in \mathcal{Y}(x^{(i)})} \left(\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y) + \text{cost}(y^{(i)}, y) \right)}_{\text{cost-augmented decoding}} \right)$$



Carnegie Mellon

Training

- Standard approach is to maximize conditional likelihood:

$$\min_{\theta} \sum_{i=1}^n \left(-\theta^\top \mathbf{f}(x^{(i)}, y^{(i)}) + \log \sum_{y \in \mathcal{Y}(x^{(i)})} \exp\{\theta^\top \mathbf{f}(x^{(i)}, y)\} \right)$$

- Another approach maximizes margin (Taskar et al., 2003):

$$\min_{\theta} \sum_{i=1}^n \left(-\theta^\top \mathbf{f}(x^{(i)}, y^{(i)}) + \max_{y \in \mathcal{Y}(x^{(i)})} \left(\theta^\top \mathbf{f}(x^{(i)}, y) + \text{cost}(y^{(i)}, y) \right) \right)$$

- Softmax-margin: replace “max” with “softmax”

$$\min_{\theta} \sum_{i=1}^n \left(-\theta^\top \mathbf{f}(x^{(i)}, y^{(i)}) + \underbrace{\log \sum_{y \in \mathcal{Y}(x^{(i)})} \exp\{\theta^\top \mathbf{f}(x^{(i)}, y) + \text{cost}(y^{(i)}, y)\}}_{\text{“cost-augmented summing”}} \right)$$



Carnegie Mellon

Training

- Standard approach is to maximize conditional likelihood:

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^n \left(-\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y^{(i)}) + \log \sum_{y \in \mathcal{Y}(x^{(i)})} \exp\{\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y)\} \right)$$

- Another approach maximizes margin (Taskar et al., 2003):

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^n \left(-\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y^{(i)}) + \max_{y \in \mathcal{Y}(x^{(i)})} \left(\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y) + \text{cost}(y^{(i)}, y) \right) \right)$$

- Softmax-margin: replace “max” with “softmax”

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^n \left(-\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y^{(i)}) + \log \sum_{y \in \mathcal{Y}(x^{(i)})} \exp\{\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y) + \text{cost}(y^{(i)}, y)\} \right)$$

Sha and Saul (2006), Povey et al. (2008)



Carnegie Mellon

Properties of Softmax-Margin

- Has a probabilistic interpretation in the minimum divergence framework ([Jelinek, 1997](#))
 - Details in technical report
- Is a bound on:
 - Max-margin
 - Conditional likelihood
 - Risk



Carnegie Mellon

Properties of Softmax-Margin

- Has a probabilistic interpretation in the minimum divergence framework (Jelinek, 1997)
 - Details in technical report
- Is a bound on:
 - Max-margin (because “softmax” bounds “max”) ✓
 - Conditional likelihood
 - Risk



Carnegie Mellon

Risk?

- **Risk** is the expected value of the cost function (Smith and Eisner, 2006; Li and Eisner, 2009):

$$\min_{\theta} \sum_{i=1}^n \mathbb{E}_{p_{\theta}(\cdot|x^{(i)})} [\text{cost}(y^{(i)}, \cdot)]$$



Carnegie Mellon

Bounding Conditional Likelihood and Risk

- Softmax-margin:

$$\sum_{i=1}^n \left(-\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y^{(i)}) + \log \sum_{y \in \mathcal{Y}(x^{(i)})} \exp\{\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y) + \text{cost}(y^{(i)}, y)\} \right)$$
$$= \underbrace{\sum_{i=1}^n \left(-\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y^{(i)}) + \log Z_i \right)}_{\text{Conditional likelihood}} + \underbrace{\sum_{i=1}^n \log \mathbb{E}_{p_i}[\exp\{\text{cost}(y^{(i)}, \cdot)\}]}_{\text{Bound on risk via Jensen's inequality}}$$

Conditional likelihood

**Bound on risk via
Jensen's inequality**



Carnegie Mellon

Bounding Conditional Likelihood and Risk

■ Softmax-margin:

$$\sum_{i=1}^n \left(-\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y^{(i)}) + \log \sum_{y \in \mathcal{Y}(x^{(i)})} \exp\{\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y) + \text{cost}(y^{(i)}, y)\} \right)$$
$$= \underbrace{\sum_{i=1}^n \left(-\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y^{(i)}) + \log Z_i \right)}_{\text{Conditional likelihood}} + \underbrace{\sum_{i=1}^n \log \mathbb{E}_{p_i}[\exp\{\text{cost}(y^{(i)}, \cdot)\}]}_{\text{Bound on risk via Jensen's inequality}}$$

Conditional likelihood

**Bound on risk via
Jensen's inequality**

**Softmax-margin is a convex bound on
max-margin, conditional likelihood, and risk**



Carnegie Mellon

Bounding Conditional Likelihood and Risk

- Softmax-margin:

$$\sum_{i=1}^n \left(-\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y^{(i)}) + \log \sum_{y \in \mathcal{Y}(x^{(i)})} \exp\{\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y) + \text{cost}(y^{(i)}, y)\} \right)$$
$$= \underbrace{\sum_{i=1}^n \left(-\boldsymbol{\theta}^\top \mathbf{f}(x^{(i)}, y^{(i)}) + \log Z_i \right)}_{\text{Conditional likelihood}} + \underbrace{\sum_{i=1}^n \log \mathbb{E}_{p_i}[\exp\{\text{cost}(y^{(i)}, \cdot)\}]}_{\text{Jensen Risk Bound}}$$

Conditional likelihood

Jensen Risk Bound

Easier to optimize than risk
(cf. Li and Eisner, 2009)



Carnegie Mellon

Implementation

- Conditional likelihood \rightarrow Softmax-margin
 - If cost function factors the same way as the features, it's easy:
 - Add additional features for the cost function
 - Keep their weights fixed
 - If not, use a simpler cost function or use approximate inference



Carnegie Mellon

Experiments

- English named-entity recognition (CoNLL 2003)
- Compared softmax-margin and Jensen risk bound with five baselines:
 - Perceptron (Collins, 2002)
 - 1-best MIRA with cost-augmented decoding (Crammer et al., 2006)
 - Max-margin via subgradient descent (Ratliff et al., 2006)
 - Conditional likelihood (Lafferty et al., 2001)
 - Risk (Xiong et al., 2009)
- For risk and Jensen risk bound, initialized using output of conditional likelihood training
- Used Hamming cost for cost function



Carnegie Mellon

Results

Method	Test F_1
Perceptron	83.98*
MIRA	85.72
Max-Margin	85.28*
Conditional Likelihood	85.46*
Risk	85.59
Jensen Risk Bound	85.65
Softmax-Margin	85.84

* Indicates significance (compared with softmax-margin)

Results

Method	Test F_1
Perceptron	83.98*
MIRA	85.72
Max-Margin	85.28*
Conditional Likelihood	85.46*
Risk	85.59
Jensen Risk Bound	85.65
Softmax-Margin	85.84

**Significant
improvement with
equal training
time and
implementation
difficulty**

*** Indicates significance (compared with softmax-margin)**



Carnegie Mellon

Results

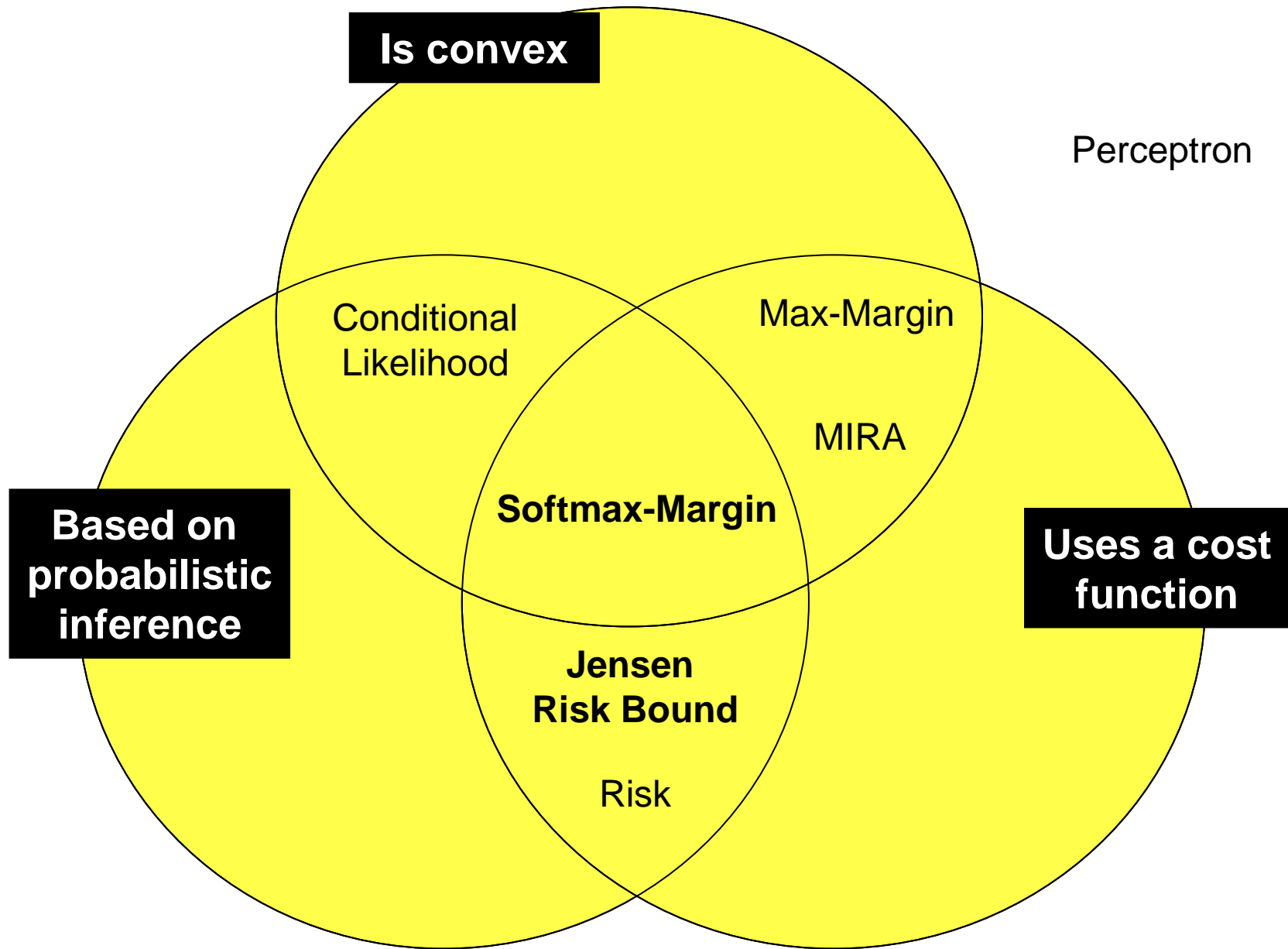
Method	Test F_1
Perceptron	83.98*
MIRA	85.72
Max-Margin	85.28*
Conditional Likelihood	85.46*
Risk	85.59
Jensen Risk Bound	85.65
Softmax-Margin	85.84

Comparable performance with half the training time

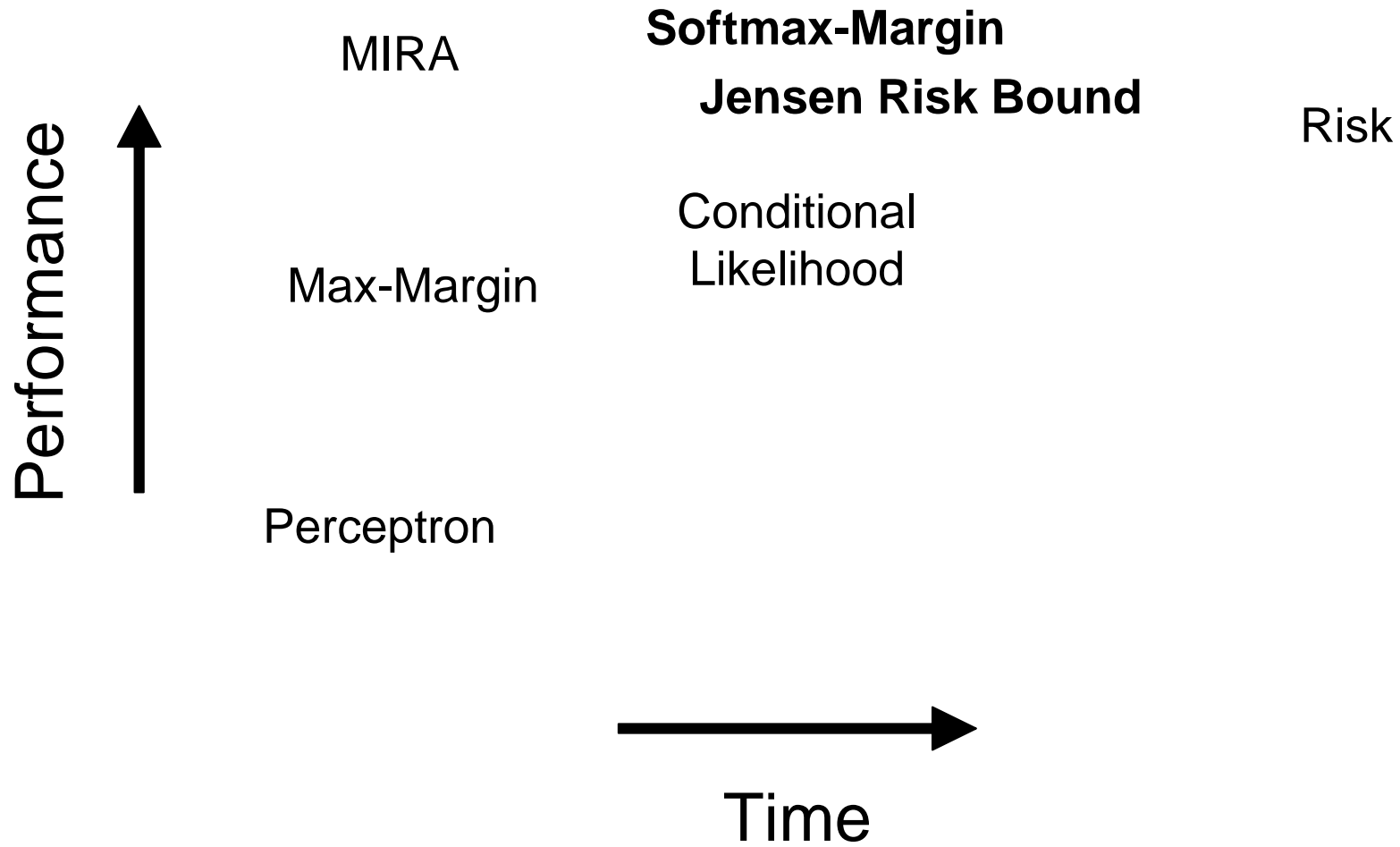
*** Indicates significance (compared with softmax-margin)**



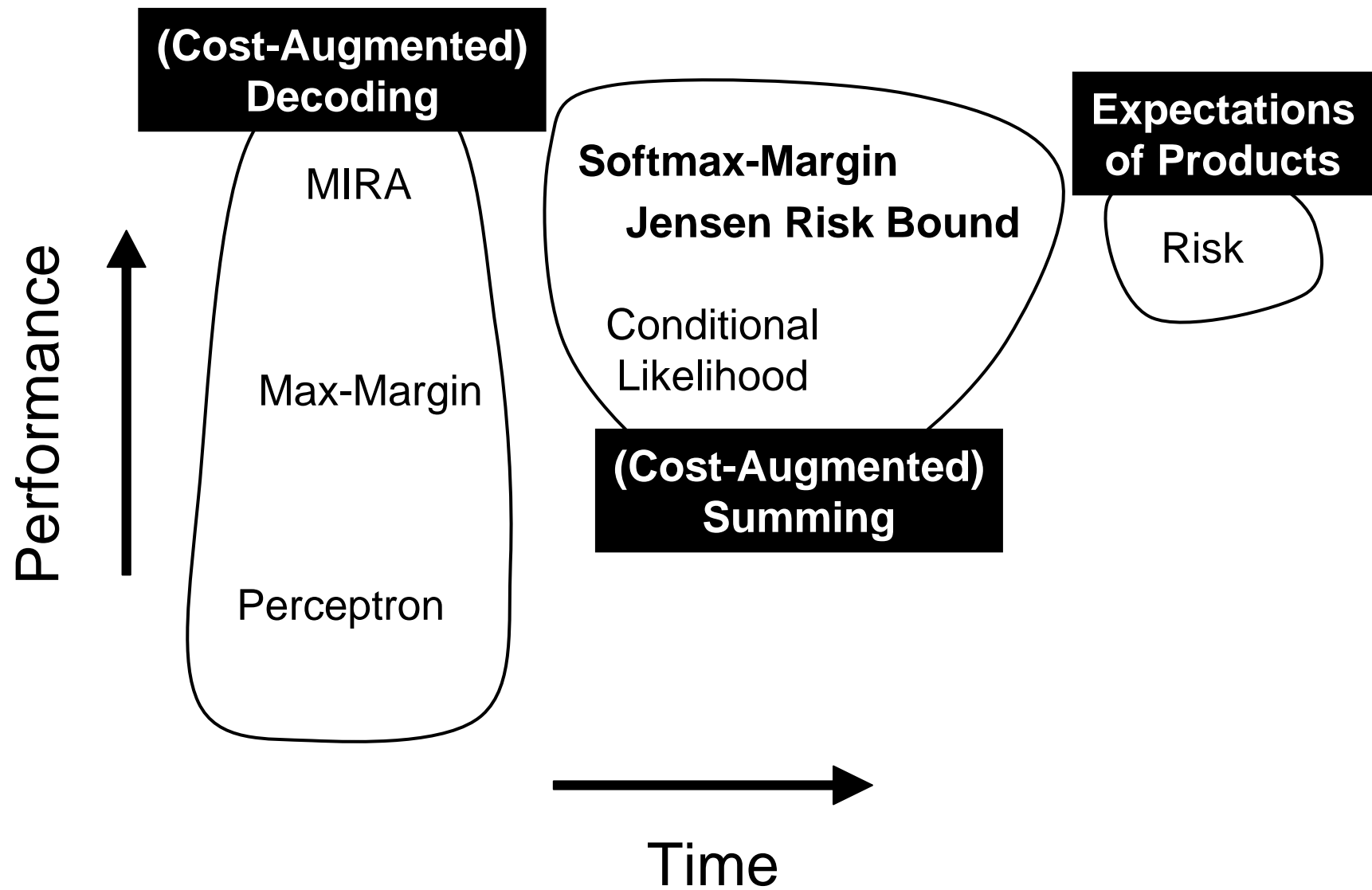
Carnegie Mellon



Carnegie Mellon



Carnegie Mellon



Carnegie Mellon

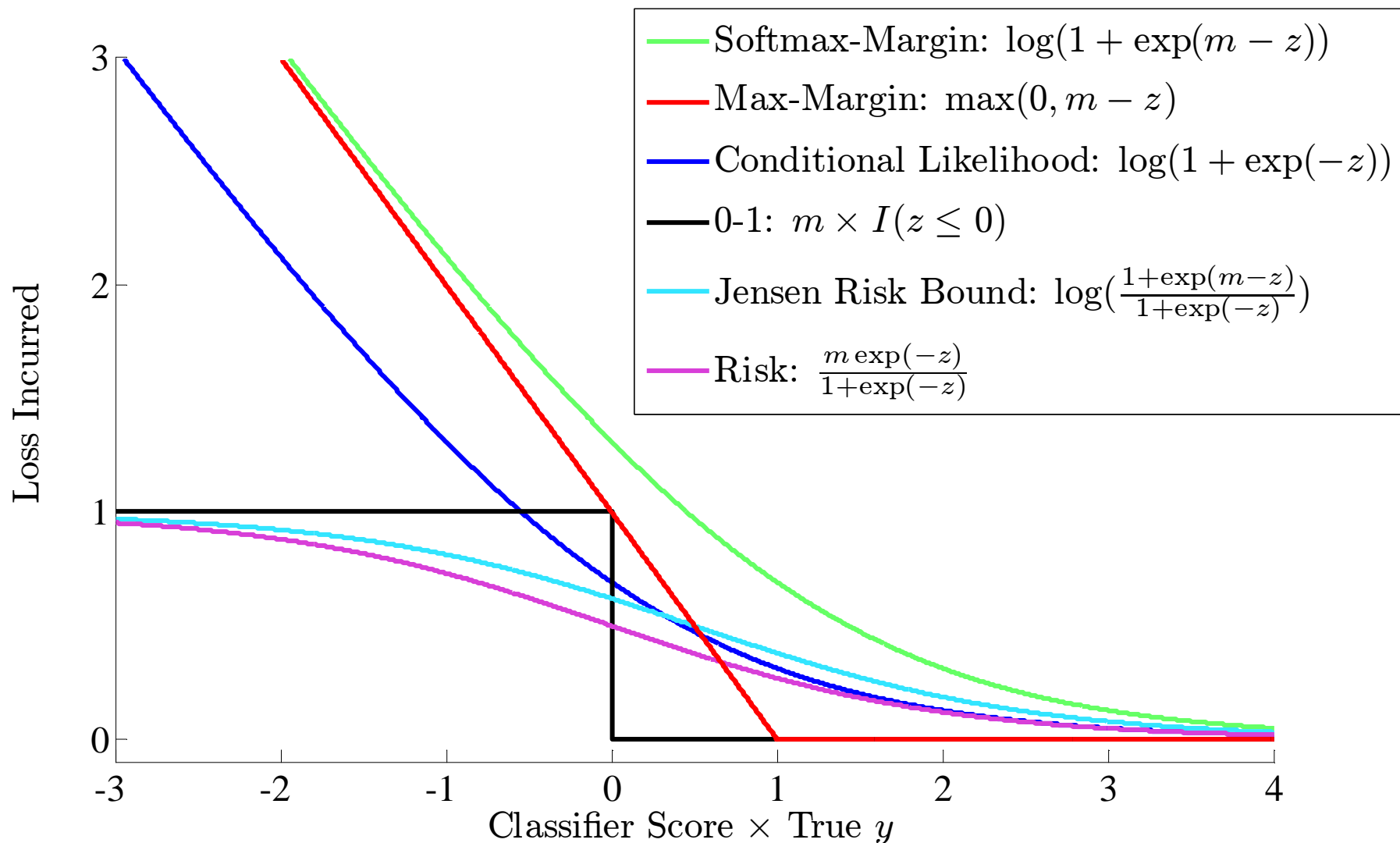
Thank you!

- See extended technical report for:
 - Probabilistic interpretation for softmax-margin in minimum divergence framework ([Jelinek, 1997](#))
 - Softmax-margin training with hidden variables
 - Additional experiments



Carnegie Mellon

Loss Functions for Binary Classification



Carnegie Mellon

Training Method	Requirements	Cost Function	Convex	Prob. Interp.
Perceptron	decoding		N/A	
MIRA	cost-augmented decoding	✓	✓	
Max-Margin	cost-augmented decoding	✓	✓	
Conditional Likelihood	summing		✓	✓
Risk	expectations of products	✓		✓
Jensen Risk Bound	cost-augmented summing	✓		✓
Softmax-Margin	cost-augmented summing	✓	✓	✓



Carnegie Mellon