

Natural Language Processing



Language Modeling III

Dan Klein – UC Berkeley



Improving on N-Grams?

- N-grams don't combine multiple sources of evidence well

$P(\text{construction} \mid \text{After the demolition was completed, the})$

- Here:
 - "the" gives syntactic constraint
 - "demolition" gives semantic constraint
 - Unlikely the interaction between these two has been densely observed in this specific n-gram
- We'd like a model that can be more statistically efficient

Maximum Entropy Models



Some Definitions

INPUTS	x_i	close the ____
CANDIDATE SET	$\mathcal{Y}(x)$	{door, table, ...}
CANDIDATES	y	table
TRUE OUTPUTS	y_i^*	door
FEATURE VECTORS	$f(x, y)$	[0 0 1 0 0 0 1 0 0 0 0 0]

$x_{-1} = \text{"the"} \wedge y = \text{"door"}$ (points to 1st '1')
 $x_{-2} = \text{"the"} \wedge y = \text{"table"}$ (points to 3rd '1')
 "close" in $x \wedge y = \text{"door"}$ (points to 7th '1')
 y occurs in x (points to 7th '1')



More Features, Less Interaction

$x = \text{closing the } ____, y = \text{doors}$

- N-Grams $x_{-1} = \text{"the"} \wedge y = \text{"doors"}$
- Skips $x_{-2} = \text{"closing"} \wedge y = \text{"doors"}$
- Lemmas $x_{-2} = \text{"close"} \wedge y = \text{"door"}$
- Caching y occurs in x



Data: Feature Impact

Features	Train Perplexity	Test Perplexity
3 gram indicators	241	350
1-3 grams	126	172
1-3 grams + skips	101	164

Exponential Form

- Weights w Features $f(x, y)$
- Linear score $w^T f(x, y)$
- Unnormalized probability

$$P(y|x, w) \propto \exp(w^T f(x, y))$$
- Probability

$$P(y|x, w) = \frac{\exp(w^T f(x, y))}{\sum_{y'} \exp(w^T f(x, y'))}$$

Likelihood Objective

- Model form:

$$P(y|x, w) = \frac{\exp(w^T f(y))}{\sum_{y'} \exp(w^T f(y'))}$$
- Likelihood of training data

$$L(w) = \log \prod_i P(y_i^* | x_i, w) = \sum_i \log \left(\frac{\exp(w^T f_i(y_i^*))}{\sum_y \exp(w^T f_i(y))} \right)$$

$$= \sum_i \left(w^T f_i(y_i^*) - \log \sum_y \exp(w^T f_i(y)) \right)$$

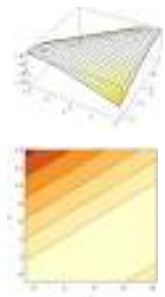
Training

History of Training

- 1990's: Specialized methods (e.g. iterative scaling)
- 2000's: General-purpose methods (e.g. conjugate gradient)
- 2010's: Online methods (e.g. stochastic gradient)

What Does LL Look Like?

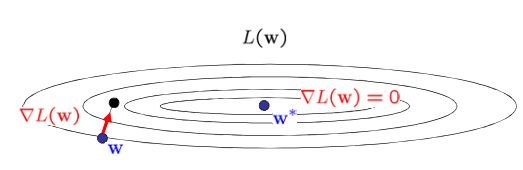
- Example
 - Data: xxxy
 - Two outcomes, x and y
 - One indicator for each
 - Likelihood



$$\log \left(\left(\frac{a^{xy}}{a^{xx} + a^{yy}} \right)^3 \cdot \frac{a^{yy}}{a^{xx} + a^{yy}} \right)$$

Convex Optimization

- The maxent objective is an unconstrained convex problem




- One optimal value*, gradients point the way


Gradients

$$L(w) = \sum_y \left(w^T R(x_i, y) - \log \sum_y \exp(w^T R(x_i, y)) \right)$$

$$\frac{\partial L(w)}{\partial w} = \sum_y \left(f(x_i, y) - \sum_y p(y|x_i) f(x_i, y) \right)$$



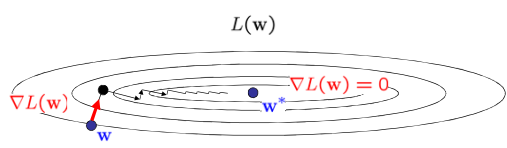
Count of features under target labels



Expected count of features under model predicted label distribution

Gradient Ascent

- The maxent objective is an unconstrained optimization problem



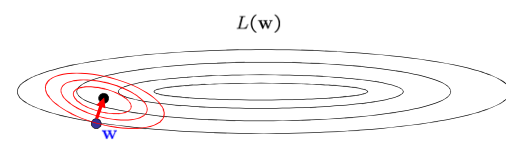
$L(w)$

$\nabla L(w)$ $\nabla L(w) = 0$ w^*

- Gradient Ascent
 - Basic idea: move uphill from current guess
 - Gradient ascent / descent follows the gradient incrementally
 - At local optimum, derivative vector is zero
 - Will converge if step sizes are small enough, but not efficient
 - All we need is to be able to evaluate the function and its derivative

(Quasi)-Newton Methods

- 2nd-Order methods: repeatedly create a quadratic approximation and solve it



$L(w)$

w



$$L(w_0) + \nabla L(w_0)^T (w - w_0) + \frac{1}{2} (w - w_0)^T \nabla^2 L(w_0) (w - w_0)$$



- E.g. LBFGS, which tracks derivative to approximate (inverse) Hessian



Regularization

Regularization Methods

- Early stopping
- L2: $LL(w) - |w|_2^2$
- L1: $LL(w) - |w|$

Regularization Effects

- Early stopping: don't do this
- L2: weights stay small but non-zero
- L1: many weights driven to zero
 - Good for sparsity
 - Usually bad for accuracy for NLP

Scaling



Why is Scaling Hard?

$$L(w) = \sum_T \left(w^T f(x_i, y_i^*) - \log \sum_Y \exp(w^T f(x_i, y)) \right)$$

- Big normalization terms
- Lots of data points



Hierarchical Prediction

- Hierarchical prediction / softmax [Mikolov et al 2013]



- Noise-Contrastive Estimation [Mnih, 2013]
- Self-Normalization [Devlin, 2014]

Image: ayende.com



Stochastic Gradient

- View the gradient as an average over data points

$$\frac{\partial L(w)}{\partial w} = \frac{1}{N} \sum_T \left(f(x_i, y_i^*) - \sum_Y P(y|x_i) f(x_i, y) \right)$$

- Stochastic gradient: take a step each example (or mini-batch)

$$\frac{\partial L(w)}{\partial w} \approx \frac{1}{1} \left(f(x_i, y_i^*) - \sum_Y P(y|x_i) f(x_i, y) \right)$$

- Substantial improvements exist, e.g. AdaGrad (Duchi, 11)

Other Methods



Neural Net LMs

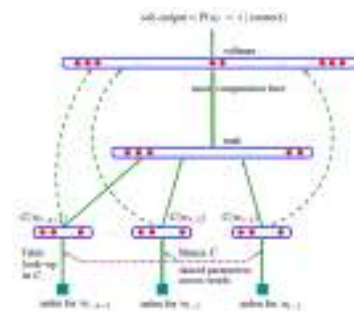


Image: (Bengio et al, 03)



Neural vs Maxent

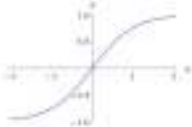
- Maxent LM

$$P(y|x, w) \propto \exp(w^T f(x, y))$$

- Neural Net LM

$$P(y|x, w) \propto \exp(B\sigma(Af(x)))$$

σ nonlinear, e.g. tanh

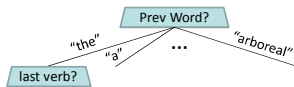


Mixed Interpolation

- But can't we just interpolate:
 - $P(w | \text{most recent words})$
 - $P(w | \text{skip contexts})$
 - $P(w | \text{caching})$
 - ...
- Yes, and people do (well, did)
 - But additive combination tends to flatten distributions, not zero out candidates



Decision Trees / Forests



- Decision trees?
 - Good for non-linear decision problems
 - Random forests can improve further [Xu and Jelinek, 2004]
 - Paths to leaves basically learn conjunctions
 - General contrast between DTs and linear models