# Algorithms for NLP

## Language Modeling III

Taylor Berg-Kirkpatrick – CMU

Slides: Dan Klein – UC Berkeley

# Announcements

- Office hours on website
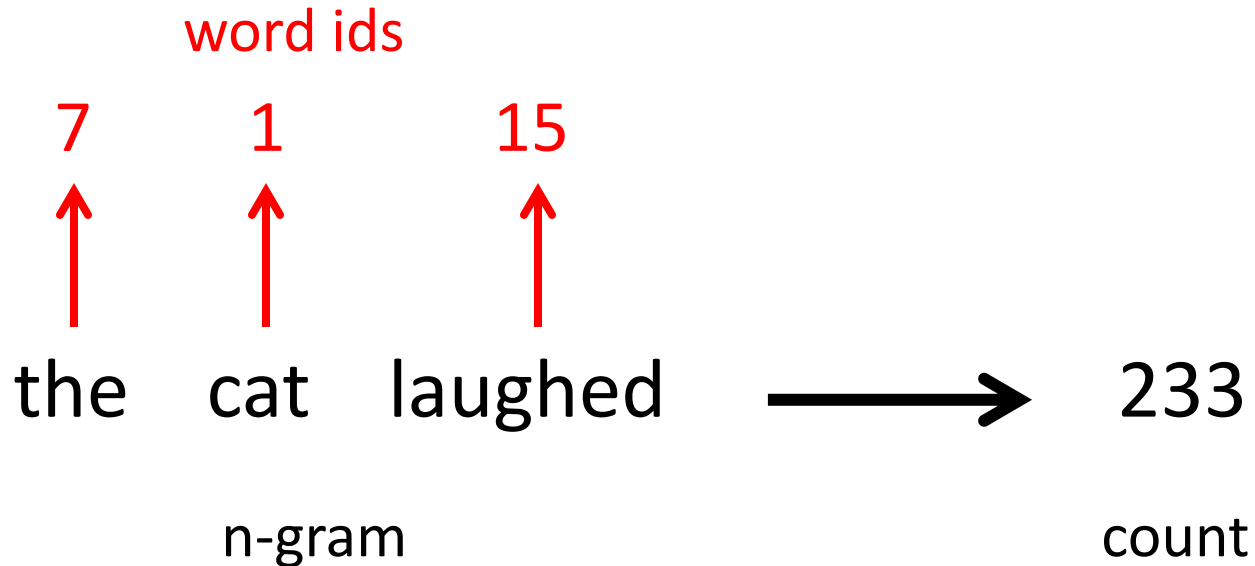  - but no OH for Taylor until next week.

# Efficient Hashing

- **Closed address hashing**
  - Resolve collisions with chains
  - Easier to understand but bigger

- **Open address hashing**
  - Resolve collisions with probe sequences
  - Smaller but easy to mess up

- **Direct-address hashing**
  - No collision resolution
  - Just eject previous entries
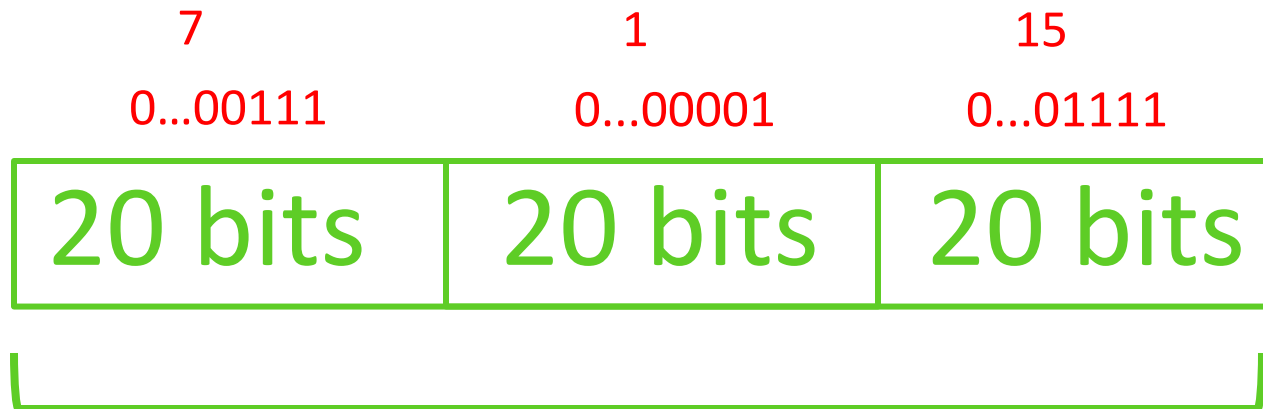  - Not suitable for core LM storage

# Integer Encodings

7    1    15

↑    ↑    ↑

the   cat   laughed   ⟶   233

n-gram             count

# Bit Packing

Got 3 numbers under $2^{20}$ to store?

| 7 | 1 | 15 |
|---|---|---|
| 0...00111 | 0...00001 | 0...01111 |
| 20 bits | 20 bits | 20 bits |

Fits in a primitive 64-bit long

# Integer Encodings

n-gram encoding

15176595 = | 20 bits | 20 bits | 20 bits |

~~the cat laughed~~ ⟶ 233

n-gram          count

# Rank Values

$$c(\text{the}) = 23135851162 < 2^{35}$$

35 bits to represent integers between 0 and $2^{35}$

60 bits                 35 bits

15176595 $\longrightarrow$ 233

n-gram encoding         count

# Rank Values

# unique counts = 770000 < $2^{20}$

20 bits to represent ranks of all counts

| rank | freq |
|------|------|
| 0 | 1 |
| 1 | 2 |
| 2 | 51 |
| 3 | 233 |

60 bits

15176595

n-gram encoding

20 bits

3

rank

# So Far

## Word indexer

| word | id |
|------|-----|
| cat | 0 |
| the | 1 |
| was | 2 |
| ran | 3 |

## Rank lookup

| rank | freq |
|------|------|
| 0 | 1 |
| 1 | 2 |
| 2 | 51 |
| 3 | 233 |

## N-gram encoding scheme

unigram:  $f(id) = id$

bigram:  $f(id_1, id_2) = \;?$

trigram:  $f(id_1, id_2, id_3) = \;?$

## Count DB

| unigram | | bigram | | trigram | |
|---------|------|---------|------|---------|------|
| 16078820 | 0381 | 16078820 | 0381 | 16078820 | 0381 |
| 15176595 | 0051 | 15176595 | 0051 | 15176595 | 0051 |
| 15176583 | 0076 | 15176583 | 0076 | 15176583 | 0076 |
| — | — | — | — | — | — |
| 16576628 | 0021 | 16576628 | 0021 | 16576628 | 0021 |
| — | — | — | — | — | — |
| 15176600 | 0018 | 15176600 | 0018 | 15176600 | 0018 |
| 16089320 | 0171 | 16089320 | 0171 | 16089320 | 0171 |
| 15176583 | 0039 | 15176583 | 0039 | 15176583 | 0039 |
| 14980420 | 0030 | 14980420 | 0030 | 14980420 | 0030 |
| — | — | — | — | — | — |
| 15020330 | 0482 | 15020330 | 0482 | 15020330 | 0482 |

# Hashing vs Sorting

query: 15176595

## Sorting

| c | val |
|---|---|
| 15176583 | 0076 |
| 15176595 | 0051 |
| 15176600 | 0018 |
| 16078820 | 0381 |
| 16089320 | 0171 |
| 16576628 | 0021 |
| 16980420 | 0030 |
| 17020330 | 0482 |
| 17176583 | 0039 |

## Hashing

| c | val |
|---|---|
| 16078820 | 0381 |
| 15176595 | 0051 |
| 15176583 | 0076 |
| — | — |
| 16576628 | 0021 |
| — | — |
| 15176600 | 0018 |
| 16089320 | 0171 |
| 15176583 | 0039 |
| 14980420 | 0030 |
| — | — |
| 15020330 | 0482 |

# Context Tries

# Tries

# Context Encodings



| 40-bits | 24-bits | |
|---|---|---|
| 548029639 | 678 | 431 |

| | 4-gram | -8.7 |
|---|---|---|

Google N-grams
- 10.5 bytes/n-gram
- 37 GB total

[Many details from Pauls and Klein, 2011]

# Context Encodings

## 1-grams

| w | val | |
|---|---|---|
| 675 | 0127 | "this" |
| 676 | 9008 | |
| 677 | 0137 | |
| 678 | 0090 | "a" |
| 679 | 1192 | |
| 680 | 0050 | "the" |
| 681 | 0040 | |
| 682 | 0201 | "is" |
| 683 | 3010 | "was" |

20 bits

## 2-grams

| | c | w | val | |
|---|---|---|---|---|
| 15176582 | 00000480 | 682 | 0065 | |
| 15176583 | 00000675 | 682 | 0808 | |
| 15176584 | 00000802 | 682 | 0012 | |
| 15176585 | 00001321 | 682 | 0400 | "is" |
| 15176586 | 00002482 | 682 | 0030 | |
| 15176587 | 00002588 | 682 | 0260 | |
| 15176588 | 00000390 | 683 | 0013 | |
| 15176589 | 00000676 | 683 | 0025 | "was" |
| 15176590 | 00000984 | 683 | 0086 | |

← 64 bits → | 20 bits

## 3-grams

| | c | w | val | |
|---|---|---|---|---|
| 42276773 | 15176583 | 678 | 0076 | |
| 42276774 | 15176595 | 678 | 0051 | |
| 42276775 | 15176600 | 678 | 0018 | |
| 42276776 | 16078820 | 678 | 0381 | "a" |
| 42276777 | 16089320 | 678 | 0171 | |
| 42276778 | 16576628 | 678 | 0021 | |
| 42276779 | 14980420 | 680 | 0030 | |
| 42276780 | 15020330 | 680 | 0482 | "the" |
| 42276781 | 15176583 | 680 | 0039 | |

← 64 bits → | 20 bits

# Compression

# Idea: Differential Compression

| c | w | val |
|---|---|---|
| 15176585 | 678 | 3 |
| 15176587 | 678 | 2 |
| 15176593 | 678 | 1 |
| 15176613 | 678 | 8 |
| 15179801 | 678 | 1 |
| 15176585 | 680 | 298 |
| 15176589 | 680 | 1 |

| $\Delta c$ | $\Delta w$ | val |
|---|---|---|
| 15176583 | 678 | 3 |
| +2 | +0 | 2 |
| +6 | +0 | 1 |
| +40 | +0 | 8 |
| +188 | +0 | 1 |
| 15176585 | +2 | 298 |
| +4 | +0 | 1 |

| $|\Delta w|$ | $|\Delta c|$ | $|val|$ |
|---|---|---|
| 40 | 24 | 3 |
| 3 | 2 | 3 |
| 3 | 2 | 3 |
| 9 | 2 | 6 |
| 12 | 2 | 3 |
| 36 | 4 | 15 |
| 6 | 2 | 3 |

| 15176585 | 678 | 563097887 | 956 | 3 | 0 | +2 | +0 | 2 | +6 | +0 | 1 | +40 | +2 | 8 | $\cdot$ $\cdot$ $\cdot$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# Variable Length Encodings

## Encoding "9"

$$\underbrace{000}_{\text{Length in Unary}} \underbrace{1001}_{\text{Number in Binary}}$$

Length in Unary

Number in Binary

Google N-grams
- 2.9 bytes/n-gram
- 10 GB total

[Elias, 75]

# Speed-Ups

# Context Encodings

## 1-grams

| w | val | |
|---|---|---|
| 675 | 0127 | "this" |
| 676 | 9008 | |
| 677 | 0137 | |
| 678 | 0090 | "a" |
| 679 | 1192 | |
| 680 | 0050 | "the" |
| 681 | 0040 | |
| 682 | 0201 | "is" |
| 683 | 3010 | "was" |

20 bits

## 2-grams

| | c | w | val | |
|---|---|---|---|---|
| 15176582 | 00000480 | 682 | 0065 | |
| 15176583 | 00000675 | 682 | 0808 | |
| 15176584 | 00000802 | 682 | 0012 | |
| 15176585 | 00001321 | 682 | 0400 | "is" |
| 15176586 | 00002482 | 682 | 0030 | |
| 15176587 | 00002588 | 682 | 0260 | |
| 15176588 | 00000390 | 683 | 0013 | |
| 15176589 | 00000676 | 683 | 0025 | "was" |
| 15176590 | 00000984 | 683 | 0086 | |

64 bits   20 bits

## 3-grams

| | c | w | val | |
|---|---|---|---|---|
| 42276773 | 15176583 | 678 | 0076 | |
| 42276774 | 15176595 | 678 | 0051 | |
| 42276775 | 15176600 | 678 | 0018 | |
| 42276776 | 16078820 | 678 | 0381 | "a" |
| 42276777 | 16089320 | 678 | 0171 | |
| 42276778 | 16576628 | 678 | 0021 | |
| 42276779 | 14980420 | 680 | 0030 | |
| 42276780 | 15020330 | 680 | 0482 | "the" |
| 42276781 | 15176583 | 680 | 0039 | |

64 bits   20 bits

# Rolling Queries

| c | w | val | suffix |
|---|---|---|---|
| 15176583 | 682 | 0065 | 00000480 |
| 15176595 | 682 | 0808 | 00000675 |
| 15176600 | 682 | 0012 | 00000802 |
| 16078820 | 682 | 0400 | 00001321 |

this is    +   a    4-gram

12438010        0045        4820

12438010 0045 ──────→ LM ──val──→ -7.8
  this is      a

         suffix

15176583 4820 ──────→ LM ──val──→ -5.4
  is a    4-gram

              suffix ──────→ 14986731

# Idea: Fast Caching

| | n-gram | probability |
|---|---|---|
| 0 | 124 80 42 1243 | -7.034 |
| 1 | 37 2435 243 21 | -2.394 |
| 2 | 804 42 4298 43 | -8.008 |

hash( 124 80 42 1243 ) =0

hash( 1423 43 42 400 ) =1

LM can be more than 10x faster w/ direct-address caching

# Approximate LMs

- Simplest option: hash-and-hope
  - Array of size K ~ N
  - (optional) store hash of keys
  - Store values in direct-address
  - Collisions: store the max
  - What kind of errors can there be?

- More complex options, like bloom filters (originally for membership, but see Talbot and Osborne 07), perfect hashing, etc

# Maximum Entropy Models

# Improving on N-Grams?

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

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

- 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

# Some Definitions

INPUTS $\mathbf{x}_i$ *close the _____*

CANDIDATE SET $\mathcal{Y}(\mathbf{x})$ *{door, table, ...}*

CANDIDATES $\mathbf{y}$ *table*

TRUE OUTPUTS $\mathbf{y}_i^*$ *door*

FEATURE VECTORS $\mathbf{f}(\mathbf{x}, \mathbf{y})$ $[0\ 0\ 1\ 0\ 0\ 0\ 1\ 0\ 0\ 0\ 0\ 0]$

$x_{-1}$="the" $\wedge$ y="door"

$x_{-1}$="the" $\wedge$ y="table"

*"close" in x* $\wedge$ y="door"

*y occurs in x*

# More Features, Less Interaction

*x = closing the _____, y = doors*

- N-Grams   $x_{-1}$*="the"* $\wedge$ y="doors"

- Skips   $x_{-2}$*="closing"* $\wedge$ y="doors"

- Lemmas   $x_{-2}$*="close"* $\wedge$ y="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 $\mathbf{w}$     Features $\mathbf{f}(\mathbf{x}, \mathbf{y})$

- Linear score $\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y})$

- Unnormalized probability

$$P(\mathbf{y}|\mathbf{x}, \mathbf{w}) \propto \exp(\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}))$$

- Probability

$$P(\mathbf{y}|\mathbf{x}, \mathbf{w}) = \frac{\exp(\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}))}{\sum_{\mathbf{y}'} \exp(\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}'))}$$

# Likelihood Objective

- Model form:

$$P(y|x, w) = \frac{\exp(w^\top f(x, y))}{\sum_{y'} \exp(w^\top f(x, y'))}$$

- Log-likelihood of training data

$$L(w) = \log \prod_i P(y_i^*|x_i, w) = \sum_i \log \left( \frac{\exp(w^\top f(x_i, y_i^*))}{\sum_{y'} \exp(w^\top f(x_i, y'))} \right)$$

$$= \sum_i \left( w^\top f(x_i, y_i^*) - \log \sum_{y'} \exp(w^\top f(x_i, y')) \right)$$

# 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

$$\log\left(\left(\frac{e^x}{e^x + e^y}\right)^3 \times \frac{e^y}{e^x + e^y}\right)$$

# Convex Optimization

- The maxent objective is an unconstrained convex problem

$$L(\mathbf{w})$$



$\nabla L(\mathbf{w})$

$\nabla L(\mathbf{w}) = 0$

$\mathbf{w}^*$

$\mathbf{w}$

- One optimal value*, gradients point the way

# Gradients

$$L(\mathbf{w}) = \sum_i \left( \mathbf{w}^\top \mathbf{f}(\mathbf{x}_i, \mathbf{y}_i^*) - \log \sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}(\mathbf{x}_i, \mathbf{y})) \right)$$

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \sum_i \left( \mathbf{f}(\mathbf{x}_i, \mathbf{y}_i^*) - \sum_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}_i) \mathbf{f}(\mathbf{x}_i, \mathbf{y}) \right)$$

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

$$L(\mathbf{w})$$



$$\nabla L(\mathbf{w})$$

$$\nabla L(\mathbf{w}) = 0$$

$$\mathbf{w}^*$$

$$\mathbf{w}$$

- 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

$$L(\mathbf{w})$$



$$L(\mathbf{w}_0) + \nabla L(\mathbf{w})^\top (\mathbf{w} - \mathbf{w}_0) + (\mathbf{w} - \mathbf{w}_0)^\top \nabla^2 L(\mathbf{w})(\mathbf{w} - \mathbf{w}_0)$$
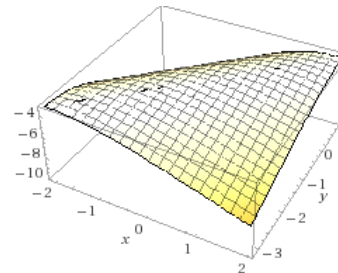
- E.g. LBFGS, which tracks derivative to approximate (inverse) Hessian
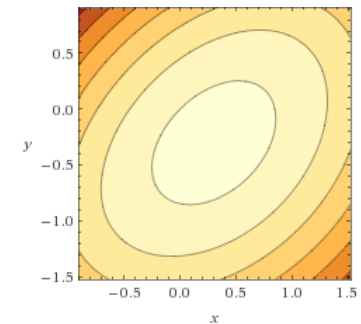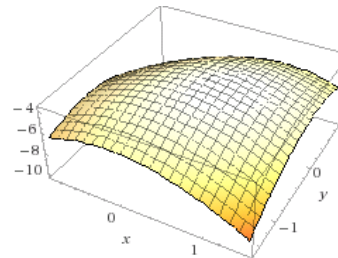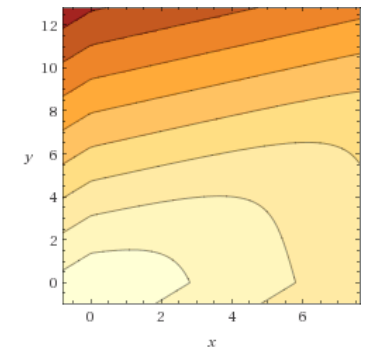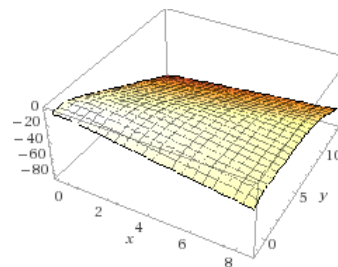
# Regularization

# Regularization Methods

- Early stopping

- L2: $L(w) - |w|_2^2$

- L1: $L(w) - |w|$

# Regularization Effects

- Early stopping: don't do this

- L2: weights stay small but non-zero

- L1: many weights driven to zero
  - Good for sparsity
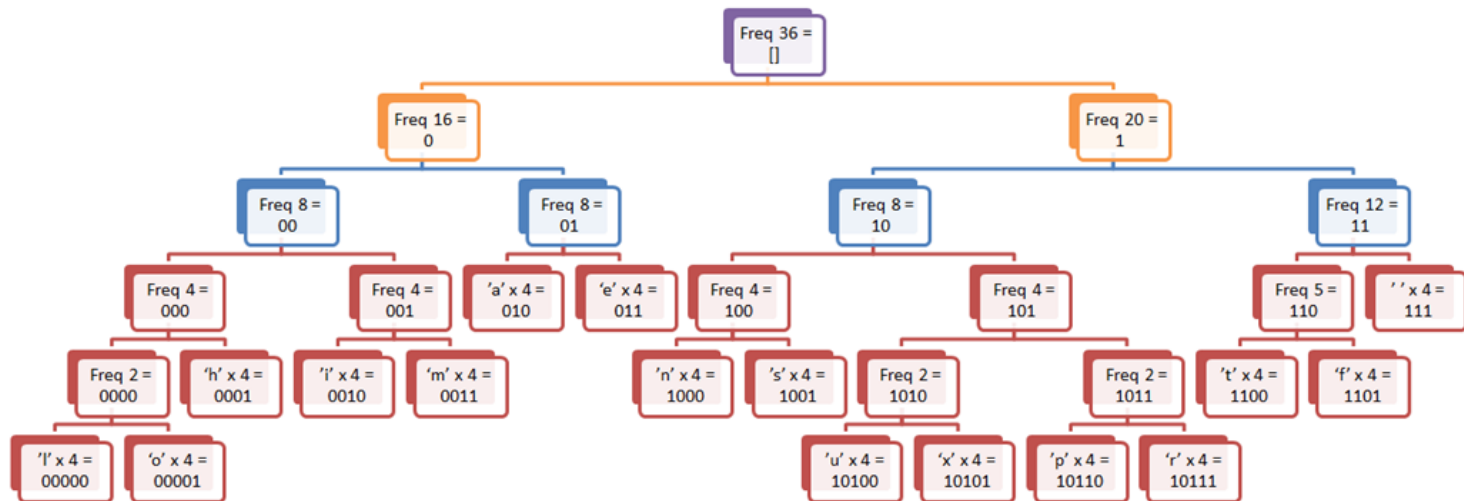  - Usually bad for accuracy for NLP

# Scaling

# Why is Scaling Hard?

$$L(\mathbf{w}) = \sum_i \left( \mathbf{w}^\top \mathbf{f}(\mathbf{x}_i, \mathbf{y}_i^*) - \log \sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}(\mathbf{x}_i, \mathbf{y})) \right)$$

- 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]

# Stochastic Gradient

- View the gradient as an average over data points

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{N} \sum_i \left( \mathbf{f}(\mathbf{x}_i, \mathbf{y}_i^*) - \sum_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}_i) \mathbf{f}(\mathbf{x}_i, \mathbf{y}) \right)$$

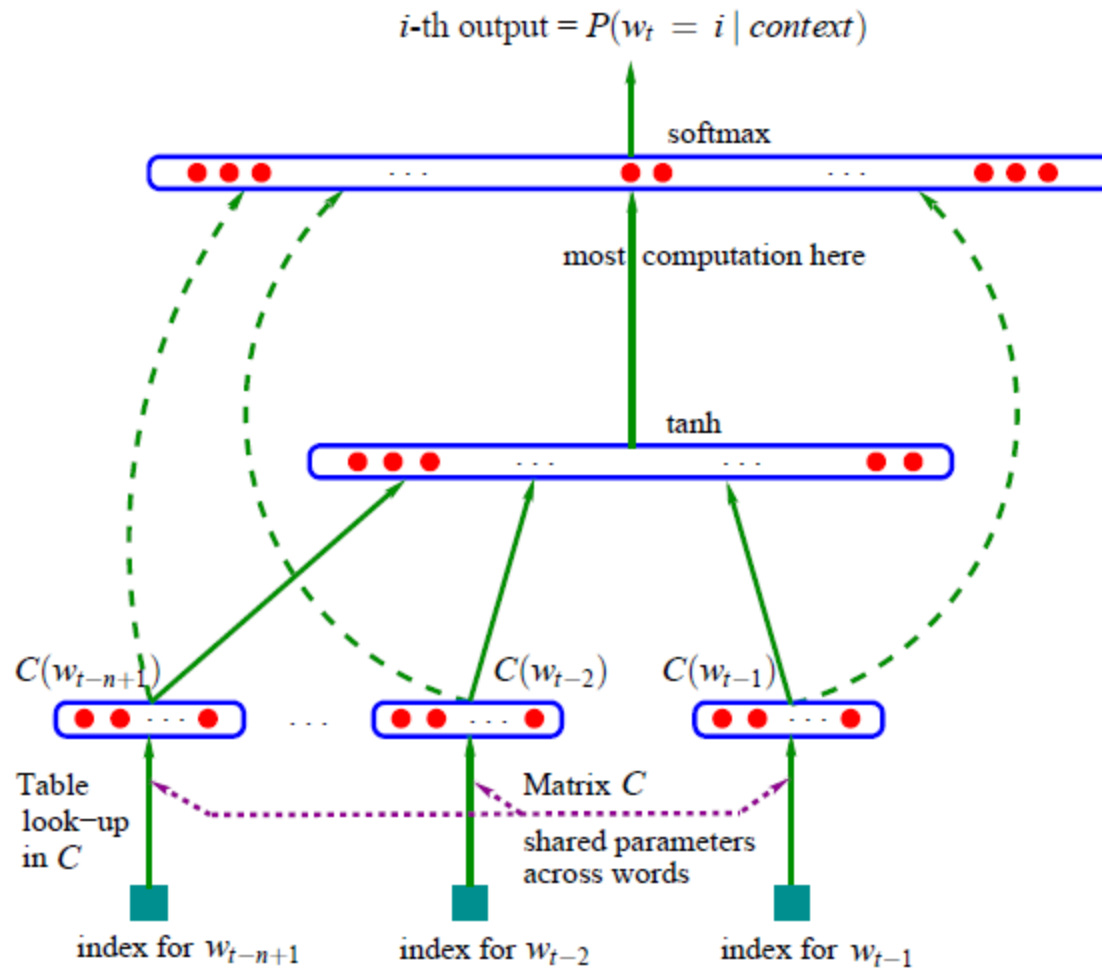- Stochastic gradient: take a step each example (or mini-batch)

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} \approx \frac{1}{1} \left( \mathbf{f}(\mathbf{x}_i, \mathbf{y}_i^*) - \sum_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}_i) \mathbf{f}(\mathbf{x}_i, \mathbf{y}) \right)$$

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

# Other Methods

# Neural Net LMs



Image: (Bengio et al, 03)

# Neural vs Maxent

- Maxent LM

$$P(\mathbf{y}|\mathbf{x}, \mathbf{w}) \propto \exp(\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}))$$

- Neural Net LM

$$P(\mathbf{y}|\mathbf{x}, \mathbf{w}) \propto \exp\left(B\sigma\left(Af(x)\right)\right)$$

$\sigma$ nonlinear, e.g. tanh

# Neural Net LMs

*man*    *door*    *doors*    *...*

$$P(y|w, x) \propto e^{Bh}$$

| $-$ | 2.3 | 1.5 | ... |
|-----|-----|-----|-----|

7.2

$$h = \sigma\left(\sum_i A_i v_i\right)$$

| 8.9 | ... |
|-----|-----|

$v_{closing}$

| 1.2 | 7.4 | ... |
|-----|-----|-----|

$v_{the}$

| $-3.3$ | 1.1 | ... |
|--------|-----|-----|

$x_{-2}$ = *closing*

$x_{-1}$= *the*

# Mixed Interpolation

- But can't we just interpolate:
    - P(w|most recent words)
    - P(w|skip contexts)
    - P(w|caching)
    - …

- Yes, and people do (well, did)
    - But additive combination tends to flatten distributions, not zero out candidates

# Decision Trees / Forests

```
                    ┌─────────────┐
                    │  Prev Word? │
                    └─────────────┘
         "the"      /    |    \      "arboreal"
              "a"  /     …    \
        ┌──────────┐
        │ last verb? │
        └──────────┘
```

- **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

- L2(0.01) 17 / 355
- L2(0.1) 27 / 172
- L2(0.5) 60 / 156
- l2(10) 296 / 265

# Maximum Entropy LMs

- Want a model over completions y given a context x:

$$P(y|x) = P(\text{ } \textit{close the } \textbf{door} \text{ } | \textit{ close the }\text{ })$$

- Want to characterize the important aspects of y = (v,x) using a feature function f

- F might include
    - Indicator of v (unigram)
    - Indicator of v, previous word (bigram)
    - Indicator whether v occurs in x (cache)
    - Indicator of v and each non-adjacent previous word
    - …