

STATE >

Recall: Multiclass Classification	Re
▶ Logistic regression: $P(y x) = \frac{\exp\left(w^{\top}f(x,y)\right)}{\sum_{y' \in \mathcal{Y}} \exp\left(w^{\top}f(x,y')\right)}$	Stochastic gradient *asc
Gradient (unregularized): $rac{\partial}{\partial w_i}\mathcal{L}(x_j,y_j^*)=f_i(x_j,y_j^*)-\mathbb{E}_y[f_i(x_j,y)]$	Adagrad:
• SVM: defined by quadratic program (minimization, so gradients are flipped) Loss-augmented decode $\xi_j = \max_{y \in \mathcal{Y}} w^\top f(x_j, y) + \ell(y, y_j^*) - w^\top f(x_j, y_j^*)$	 SGD/AdaGrad have a batches Large batches (>50 ex but bigger batches of you make fewer paraget
Subgradient (unregularized) on jth example $= f_i(x_j,y_{ m max}) - f_i(x_j,y_j^st)$	Shuffling: online metho

10 MIT

Stochastic gradient *ascent* • Adagrad: • SGD/AdaGrad have a batch size parameter • Large batches (>50 examples): can parallelize within batch • ...but bigger batches often means more epochs required because you make fewer parameter updates • Shuffling: online methods are sensitive to dataset order, shuffling helps!







What is this good for?	Sequence Models
Text-to-speech: record, lead	lnput $\mathbf{x} = (x_1,, x_n)$ Output $\mathbf{y} = (y_1,, y_n)$
Preprocessing step for syntactic parsers	
Domain-independent disambiguation for other tasks	POS tagging: x is a sequence of words, y is a sequence of tags
 (Very) shallow information extraction 	▶ Today: generative models P(<i>x, y</i>); discriminative models next time

Hidden Markov Models

• Input $\mathbf{x} = (x_1,...,x_n)$ Output $\mathbf{y} = (y_1,...,y_n)$

- Model the sequence of y as a Markov process (dynamics model)
- Markov property: future is conditionally independent of the past given the present

$$(y_1 \rightarrow y_2 \rightarrow y_3) \quad P(y_3|y_1, y_2) = P(y_3|y_2)$$

- Lots of mathematical theory about how Markov chains behave
- If y are tags, this roughly corresponds to assuming that the next tag only depends on the current tag, not anything before

Hidden Markov Models

Input $\mathbf{x} = (x_1, ..., x_n)$ Output $\mathbf{y} = (y_1, ..., y_n)$



Initial Transition Emission distribution probabilities probabilities

- Observation (x) depends only on current state (y)
- Multinomials: tag x tag transitions, tag x word emissions
- P(x|y) is a distribution over all words in the vocabulary

 not a distribution over features (but could be!)



Emissions in POS Tagging

NNP VBZ NN NNS CD NN Fed raises interest rates 0.5 percent

- Emissions $P(x \mid y)$ capture the distribution of words occurring with a given tag
- P(word | NN) = (0.05 person, 0.04 official, 0.03 interest, 0.03 percent ...)
- When you compute the posterior for a given word's tags, the distribution favors tags that are more likely to generate that word
- How should we smooth this?

3



Many neural sequence models depend on entire previous tag sequence, need to use approximations like beam search











slide credit: Vivek Srikumar









Errors	Remaining Errors
JJ NN NNP NNPS RB RP IN VB VBD VBN VBP Total JJ 0 177 56 0 61 2 5 10 15 108 0 488 NN 244 0 103 0 12 1 1 29 5 6 19 525 NNP 107 106 0 132 5 0 7 5 1 2 0 427 NNPS 1 0 110 0 0 0 0 0 0 104 2 9 5 6 19 525 RB 72 21 7 0 0 16 138 1 0 0 0 323 VB 17 64 9 0 2 0 1 0 1 221 YB5 34 3 1 1 0 2 49 6 3 0 104 12 21 166 <	 Lexicon gap (word not seen with that tag in training) 4.5% Unknown word: 4.5% Could get right: 16% (many of these involve parsing!) Difficult linguistics: 20% VBD / VBP? (past or present?) They set up absurd situations, detached from reality Underspecified / unclear, gold standard inconsistent / wrong: 58% adjective or verbal participle? JJ / VBN? a \$ 10 million fourth-quarter charge against discontinued operations Manning 2011 "Part-of-Speech Tagging from 97% to 100%: Is It Time for Some Linguistics?"

