

$$\begin{array}{c}
\end{array}$$
Recall: Multiclass Classification
$$\begin{array}{c}
\end{array}$$
Logistic regression: $P(y|x) = \frac{\exp\left(w^{\top}f(x,y)\right)}{\sum_{y'\in\mathcal{Y}}\exp\left(w^{\top}f(x,y')\right)}$
Gradient (unregularized): $\frac{\partial}{\partial w_i}\mathcal{L}(x_j,y_j^*) = f_i(x_j,y_j^*) - \mathbb{E}_y[f_i(x_j,y)]$

$$\begin{array}{c}
\end{array}$$
SVM: defined by quadratic program (minimization, so gradients are flipped)
Loss-augmented decode
$$\xi_j = \max_{y\in\mathcal{Y}}w^{\top}f(x_j,y) + \ell(y,y_j^*) - w^{\top}f(x_j,y_j^*)$$
Subgradient (unregularized) on *j*th example $= f_i(x_j, y_{\max}) - f_i(x_j, y_j^*)$

$$\begin{array}{c}
\end{array}$$

Optimization

 $w \leftarrow$

Stochastic gradient *ascent*

$$w + \alpha g, \quad g = \frac{\partial}{\partial w} \mathcal{L}$$

- Very simple to code up
- "First-order" technique: only relies on having gradient
- Difficult to tune step size
- Newton's method

- Second-order technique
- Optimizes quadratic instantly
- $w \leftarrow w + \left(\frac{\partial^2}{\partial w^2} \mathcal{L}\right)^{-1} g$ Inverse Hessian: *n* x *n* mat, expensive!
- Quasi-Newton methods: L-BFGS, etc.
- Approximate inverse Hessian with gradients over time

AdaGrad

Other techniques for optimizing deep models — more later!

Duchi et al. (2011)



((;;;))





What is this good for? Text-to-speech: record, lead Preprocessing step for syntactic parsers Domain-independent disambiguation for other tasks (Very) shallow information extraction Tod

Sequence Models

- ightarrow Input $\mathbf{x} = (x_1, ..., x_n)$ Output $\mathbf{y} = (y_1, ..., y_n)$
- POS tagging: x is a sequence of words, y is a sequence of tags (most of the time...)
- > Today: generative models P(x, y); discriminative models next time





Transitions in POS Tagging

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

Should y be a single tag?

- Trigram model: y₁ = (<S>, NNP), y₂ = (NNP, VBZ), ...
- P((VBZ, NN) | (NNP, VBZ)) more context! Noun-verb-noun S-V-O
- Tradeoff between model capacity and data size















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 RB 72 21 7 0 0 16 138 1 0 0 323 VB 17 64 9 0 2 0 1 443 2 166 VBN 101 3 3 0 0 0 3 108 0 104 </td <td> 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?" </td>	 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?"

