













Multiclass Classification

• Not all classes may even be separable using this approach



۲

 Can separate 1 from 2+3 and 2 from 1+3 but not 3 from the others (with these features)



- \blacktriangleright Formally: instead of two labels, we have an output space ${\mathcal Y}$ containing a number of possible classes
- Same machinery that we'll use later for exponentially large output spaces, including sequences and trees features depend on choice

of label now! note: this

isn't the gold label

▶ Decision rule: $\operatorname{argmax}_{u \in \mathcal{V}} w^{\top} f(x, y)$ ←

۲

- Multiple feature vectors, one weight vector
- Can also have one weight vector per class: $\operatorname{argmax}_{y \in \mathcal{Y}} w_y^\top f(x)$

Different Weights vs. Different Features

- Different features: $\operatorname{argmax}_{y \in \mathcal{Y}} w^{\top} f(x, y)$
 - Suppose \mathcal{Y} is a structured label space (part-of-speech tags for each word in a sentence). f(x,y) extracts features over shared parts of these
- Different weights: $\operatorname{argmax}_{y \in \mathcal{Y}} w_y^\top f(x)$
 - ▶ Generalizes to neural networks: f(x) is the first n-1 layers of the network, then you multiply by a final linear layer at the end
- For linear multiclass classification with discrete classes, these are identical









| Training | Logistic Regression: Summary |
|--|--|
| $\begin{aligned} \frac{\partial}{\partial w_i} \mathcal{L}(x_j, y_j^*) &= f_i(x_j, y_j^*) - \sum_y f_i(x_j, y) P_w(y x_j) \\ \text{too many drug trials, too few patients} & y^* = \text{Health} \\ f(x, y = \text{Health}) &= [1, 1, 0, 0, 0, 0, 0, 0, 0] \\ f(x, y = \text{Sports}) &= [0, 0, 0, 1, 1, 0, 0, 0, 0] \\ \text{gradient:} \\ [1, 1, 0, 0, 0, 0, 0, 0, 0] &- 0.2 [1, 1, 0, 0, 0, 0, 0, 0, 0] \\ &- 0.3 [0, 0, 0, 0, 0, 1, 1, 0] \\ &= [0.8, 0.8, 0, -0.5, -0.5, 0, -0.3, -0.3, 0] \end{aligned}$ | Model: P_w(y x) = exp (w[⊤]f(x, y)) ∑_{y'∈𝔅} exp (w[⊤]f(x, y')) Inference: argmax_yP_w(y x) Learning: gradient ascent on the discriminative log-likelihood f(x, y[*]) - E_y[f(x, y)] = f(x, y[*]) - ∑_y[P_w(y x)f(x, y)] "towards gold feature value, away from expectation of feature value" |





Revisiting Generative vs. Discriminative Models

Learning in Probabilistic Models

- ▶ So far we have talked about discriminative classifiers (e.g., logistic regression which models P(y|x))
- Cannot analytically compute optimal weights for such models, need to use gradient descent
- What about generative models?





Maximum Likelihood Estimation

- ▶ Data points (x_i, y_i) provided (*j* indexes over examples)
- \blacktriangleright Find values of $P(y), \ P(x_i|y)$ that maximize data likelihood (generative):

$$\prod_{j=1}^{m} P(y_j, x_j) = \prod_{j=1}^{m} P(y_j) \left[\prod_{i=1}^{n} P(x_{ji}|y_j) \right]$$

data points (j) features (i) i th feature of jth example

Equivalent to maximizing logarithm of data likelihood:

$$\sum_{j=1}^{m} \log P(y_j, x_j) = \sum_{j=1}^{m} \left[\log P(y_j) + \sum_{i=1}^{n} \log P(x_{ji}|y_j) \right]$$

Can do this by counting and normalizing distributions!

| Next time: HMMs / POS tagging Locally-normalized generative models, so easy to estimate from data First thing we have that we could plausibly sample real sentences from In 2 lectures: CRFs (NER) You've now seen everything you need to implement multi-class classification models | | Summary |
|---|--|---|
| Locally-normalized generative models, so easy to estimate from data First thing we have that we could plausibly sample real sentences from In 2 lectures: CRFs (NER) You've now seen everything you need to implement multi-class classification models | • Next time: HM | Ms / POS tagging |
| First thing we have that we could plausibly sample real sentences from In 2 lectures: CRFs (NER) You've now seen everything you need to implement multi-class classification models | Locally-norr | nalized generative models, so easy to estimate from data |
| In 2 lectures: CRFs (NER) You've now seen everything you need to implement multi-class classification models | First thing w | e have that we could plausibly sample real sentences from |
| You've now seen everything you need to implement multi-class classification models | In 2 lectures: C | RFs (NER) |
| | You've now seen everything you need to implement multi-class classification models | |
| | | |
| | | |