# Multiclass & Structured Prediction

DS-GA 1003 Machine Learning

CDS, NYU

March 31, 2021

# Multiclass Hypothesis Space: Reframed

- **General [Discrete] Output Space:** $\mathcal{Y}$
- **Base Hypothesis Space:** $\mathcal{H} = \{h : \mathcal{X} \times \mathcal{Y} \to \mathsf{R}\}$
    - $h(x, y)$ gives **compatibility score** between input $x$ and output $y$
- **Multiclass Hypothesis Space**

$$\mathcal{F} = \left\{ x \mapsto \underset{y \in \mathcal{Y}}{\arg\max}\, h(x, y) \mid h \in \mathcal{H} \right\}$$

- Final prediction function is an $f \in \mathcal{F}$.
- For each $f \in \mathcal{F}$ there is an underlying compatibility score function $h \in \mathcal{H}$.

# Part-of-speech (POS) Tagging

- Given a sentence, give a part of speech tag for each word:

| $x$ | [START] | He | eats | apples |
|---|---|---|---|---|
| | $x_0$ | $x_1$ | $x_2$ | $x_3$ |
| $y$ | [START] | Pronoun | Verb | Noun |
| | $y_0$ | $y_1$ | $y_2$ | $y_3$ |

- $\mathcal{V} = \{$all English words$\} \cup \{$[START],"."$\}$
- $\mathcal{P} = \{$START, Pronoun, Verb, Noun, Adjective$\}$
- $\mathcal{X} = \mathcal{V}^n$, $n = 1, 2, 3, \ldots$ [Word sequences of any length]
- $\mathcal{Y} = \mathcal{P}^n$, $n = 1, 2, 3, \ldots$ [Part of speech sequence of any length]

# Structured Prediction

- A **structured prediction** problem is a multiclass problem in which $\mathcal{Y}$ is very large, but has (or we assume it has) a certain structure.

- For POS tagging, $\mathcal{Y}$ grows exponentially in the length of the sentence.

- Typical **structure** assumption: The POS labels form a Markov chain.
    - i.e. $y_{n+1} \mid y_n, y_{n-1}, \ldots, y_0$ is the same as $y_{n+1} \mid y_n$.

# Local Feature Functions: Type 1

- A "type 1" **local feature** only depends on
    - the label at a single position, say $y_i$ (label of the $i$th word) and
    - $x$ at any position

- Example:

$$\phi_1(i, x, y_i) = 1(x_i = \text{runs})1(y_i = \text{Verb})$$
$$\phi_2(i, x, y_i) = 1(x_i = \text{runs})1(y_i = \text{Noun})$$
$$\phi_3(i, x, y_i) = 1(x_{i-1} = \text{He})1(x_i = \text{runs})1(y_i = \text{Verb})$$

# Local Feature Functions: Type 2

- A "type 2" **local feature** only depends on
  - the labels at 2 consecutive positions: $y_{i-1}$ and $y_i$
  - $x$ at any position

- Example:

$$
\begin{aligned}
\theta_1(i, x, y_{i-1}, y_i) &= 1(y_{i-1} = \text{Pronoun})1(y_i = \text{Verb}) \\
\theta_2(i, x, y_{i-1}, y_i) &= 1(y_{i-1} = \text{Pronoun})1(y_i = \text{Noun})
\end{aligned}
$$

# Local Feature Vector and Compatibility Score

- At each position $i$ in sequence, define the **local feature vector**:

$$\Psi_i(x, y_{i-1}, y_i) = (\phi_1(i, x, y_i), \phi_2(i, x, y_i), \ldots,$$
$$\theta_1(i, x, y_{i-1}, y_i), \theta_2(i, x, y_{i-1}, y_i), \ldots)$$

- **Local compatibility score** for $(x, y)$ at position $i$ is $\langle w, \Psi_i(x, y_{i-1}, y_i) \rangle$.

# Sequence Compatibility Score

- The **compatibility score** for the pair of sequences $(x, y)$ is the sum of the local compatibility scores:

$$\sum_i \langle w, \Psi_i(x, y_{i-1}, y_i) \rangle$$

$$= \left\langle w, \sum_i \Psi_i(x, y_{i-1}, y_i) \right\rangle$$

$$= \langle w, \Psi(x, y) \rangle,$$

where we define the sequence feature vector by

$$\Psi(x, y) = \sum_i \Psi_i(x, y_{i-1}, y_i).$$

- So we see this is a special case of linear multiclass prediction.

# Sequence Target Loss

- How do we assess the loss for prediction sequence $y'$ for example $(x, y)$?

- **Hamming loss** is common:

$$\Delta(y, y') = \frac{1}{|y|} \sum_{i=1}^{|y|} 1(y_i \neq y_i')$$

- Could generalize this as

$$\Delta(y, y') = \frac{1}{|y|} \sum_{i=1}^{|y|} \delta(y_i, y_i')$$

## What remains to be done?

- To compute predictions, we need to find

$$\arg\max_{y \in \mathcal{Y}} \langle w, \Psi(x, y) \rangle.$$

- This is straightforward for $|\mathcal{Y}|$ small.

- Now $|\mathcal{Y}|$ is exponentially large.

- Because $\Psi$ breaks down into local functions only depending on 2 adjacent labels,
  - we can solve this efficiently using dynamic programming.
  - (Similar to Viterbi decoding.)

- Learning can be done with SGD and a similar dynamic program.

# References

- DS-GA 1003 Machine Learning Spring 2019