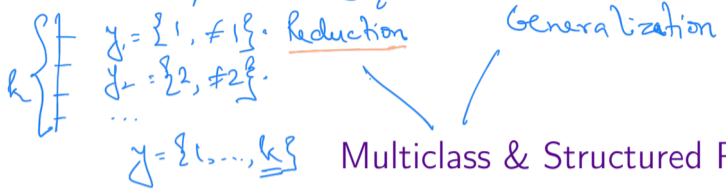


(K) One vs All/Rest $\rightarrow \operatorname{argmax}_i h_i(x, y)$

$k(k-1)/2$ AvA (K2)



$\underline{\underline{h_k(x)}}$
 $h_j(x) - h_i(x)$

Multiclass & Structured Prediction

Binary classification

DS-GA 1003 Machine Learning

Prog

Cons

- Perceptron
- SVM

$\{0, 1\}$
 $\{-1, 1\}$

CDS, NYU

- corrections
- no gradient

not linearly separable

March 31, 2021

- Margin based
- sparse solution in span of support vectors

no probabilities

- Logistic Regression

Multinomial

- probabilities
- Max likelihood Estimation

[K] Annotated

Multiclass Hypothesis Space: Reframed

- General [Discrete] Output Space: $y = \{1, \dots, k\}$
- Base Hypothesis Space: $\mathcal{H} = \{h: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}\}$
 - $h(x, y)$ gives **compatibility score** between input x and output y
- Multiclass Hypothesis Space *of energy*

Multiclass
Perceptron
SVM

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

- 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

Structured Prediction

Sequences \rightarrow NLP, DNA seq
Trees, graphs

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

x	<u>[START]</u> x_0	He x_1	<u>eats</u> x_2	apples x_3
y	<u>[START]</u> y_0	Pronoun y_1	Verb y_2	Noun y_3

- $\mathcal{V} = \{\text{all English words}\} \cup \{[\text{START}], ". ."]\}$ - Vocabulary
- $\mathcal{P} = \{\text{START, Pronoun, Verb, Noun, Adjective}\}$ - POS
- $\mathcal{X} = \mathcal{V}^n, n = 1, 2, 3, \dots$ [Word sequences of any length]
- $\mathcal{Y} = \mathcal{P}^n, n = 1, 2, 3, \dots$ [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. $y_{n+1} | y_n, y_{n-1}, \dots, y_0$ is the same as $y_{n+1} | y_n$.

Markov Assumption

Local Feature Functions: Type 1 — Unary

- 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: $\langle s \rangle$ He eats apples $\langle s \rangle$ He runs fast

$$\phi_1(i, x, y_i) = 1(x_i = \text{runs})1(y_i = \text{Verb})$$

$$\phi_2(i, x, y_i) = 1(x_i = \text{runs})1(y_i = \text{Noun})$$

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

⋮

$$\Phi(i, x, y_i) = (\phi_1(i, x, y_i), \phi_2(i, x, y_i), \phi_3(i, x, y_i))$$

e.g. $\Phi(2, x, \text{eats}) = \left(\begin{matrix} 0 \\ 0 \end{matrix}, \begin{matrix} 0 \\ 1 \end{matrix}, \begin{matrix} 1 \\ 0 \end{matrix} \right)$

Note: In the example, 'eats' is labeled as Noun, Verb.

Local Feature Functions: Type 2 — Marker

- 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:

$$\theta_1(i, x, y_{i-1}, y_i) = 1(y_{i-1} = \text{Pronoun})1(y_i = \text{Verb})$$

$$\theta_2(i, x, y_{i-1}, y_i) = 1(y_{i-1} = \text{Pronoun})1(y_i = \text{Noun})$$

$$\theta_3(i, x, \text{Noun}) = (0, 1)$$

Local Feature Vector and Compatibility Score

- At each position i in sequence, define the **local feature vector**: (class/label dependent)

$$\Psi_i(x, y_{i-1}, y_i) = \left(\underbrace{\phi_1(i, x, y_i), \phi_2(i, x, y_i), \dots}_{\text{Type 1}}, \underbrace{\theta_1(i, x, y_{i-1}, y_i), \theta_2(i, x, y_{i-1}, y_i), \dots)}_{\text{Type 2}} \right)$$

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

$$h: \mathcal{X} \rightarrow \mathbb{R}$$

$$\langle w^T \Psi_i \rangle$$

Sequence Compatibility Score

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

$$\begin{aligned} & \frac{\sum_i \langle w, \Psi_i(x, y_{i-1}, y_i) \rangle}{1} \\ &= \left\langle w, \sum_i \Psi_i(x, y_{i-1}, y_i) \right\rangle \\ &= \langle w, \Psi(x, y) \rangle, \quad - h(x) \end{aligned}$$

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

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

δ 0-1 loss \rightarrow Perceptron
 δ Hinge loss \rightarrow SVM

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

\uparrow
 $f \in \{x\} \rightarrow \operatorname{argmax}_{f \in \mathcal{Y}} h$

$\operatorname{argmax}_{f \in \mathcal{Y}} h(x, y)$ e.g.

P	V	N
P	N	N
0	1	0

 $= 1/3$

What remains to be done?

- To compute predictions, we need to find

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

learnt

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

- Because Ψ breaks down into local functions only depending on 2 adjacent labels,

- we can solve this efficiently using dynamic programming.
- (Similar to Viterbi decoding.)

reduces time complexity from exponential to polynomial

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

- DS-GA 1003 Machine Learning Spring 2019