# Lecture 7

## Kernel Machines - Structured Prediction

$\boxed{\text{Brace}}$ $\rightarrow$ $\boxed{\text{BRACE}}$

Input

output: must be a word!

$\rightarrow$ only allow certain structured outputs

Gene Sequences:

GAGTCA

$\rightarrow$ human protein structure

$\boxed{\text{The cat sat on the mat}}$ $\rightarrow$ only allow grammatically correct sentences

Kernel Input can be structured $\rightarrow$ previous lecture

e.g., weighted degree kernels

## Kernel Based Structured Prediction:

$$k : (X \times Y) \times (X \times Y) \rightarrow \mathbb{R}$$

PSD + symm.: Feature map $\phi : X \times Y \rightarrow \mathbb{R}^h$

Model: $\quad f(y|x) = w^T \phi(x,y)$

## Multi-Class Classification

$X \in \mathbb{R}^d$, $Y = \{1, ..., C\}$

input (e.g. image)   output (e.g. "dog")

$k((x,y),(x',y')) = \langle \phi(x,y), \phi(x',y') \rangle$  $\leftarrow$ general

$K((x,y),(x',y')) = k(x,x') \mathbb{1}(y=y')$  $\leftarrow$ for multiclass

Feature Map: $\psi(x,y) = \begin{pmatrix} \phi(x) \cdot \mathbb{1}(y=1) \\ \vdots \\ \phi(x) \cdot \mathbb{1}(y=c) \end{pmatrix}$

$y|x = \underset{y}{\text{argmax}} \langle w, \psi(x,y) \rangle = \underset{c}{\text{argmax}} \, w_c^T \phi(x)$

# How to learn with a structured output model?

Model: $f(y|x) = w^T \varphi(x, y)$   Prediction $y|x = \arg\max\limits_{y \in Y} f(y|x)$

Find largest-margin model by solving

$$\min_{w} \frac{1}{2} \|w\|^2 + C \sum_{n=1}^{N} \xi_n$$

Constraints: $\forall_{n=1}^{N}, \forall_{y \in y_n}: \underbrace{w^T \varphi(x_n, y_n)}_{A} - \underbrace{w^T \varphi(x_n, y)}_{B} \geq 1 - \xi_n$

$y_n$ = ground truth   $y$ = all possible evaluation

A: Score of ground truth   B: score of all possible evaluation

A - B: difference of scores, $A - B \geq 1(-\xi_n)$

$\xi_n$: Slack term: account for wrong labels, noise, ...

$\min_{w} \frac{1}{2} \|w\|^2 \rightarrow$ maximize margin

larger $C \rightarrow$ smaller $\xi_n \rightarrow$ hard margin constraint

## Structured Prediction vs. HMM

HMM Pro? - unsupervised

    - once model learned, procedure for prediction determined

Struct. Pred. Pro?

    - feature map and loss function give more flexibility for structure of problem

    - model parameter $w \in \mathbb{R}^h$ can be actively optimized for best performance on supervised task

## Summary

- structured output learning enable prediction of structured outputs (trees, sequences)
- assign matching scores to input-output pairs
- supervised   - kernel-based framework
- difficulty: efficiently infer which output $y \in Y$ maximizes the score $f(y|x)$