

Other SSL approaches  
Bayesian Learning  
Conjugate priors

Admin

AS: due now

Project Proposal: due monday

AG: due wednesday  $\Rightarrow$  is a tutorial this week

Midterm: November 10  $\Rightarrow$  (or hand in Thurs/Friday during tutorials)

⊗ Marking  $\rightarrow$  don't be a troll!  
 $\rightarrow$  give half marks!  
 $\rightarrow$  give Scott full marks!

Cross-Validation with regularization:

full dataset:

$$\sum_{i=1}^n f(y_i; w^T x_i) + \frac{\lambda}{2} \|w\|^2$$

Cross Validation Fold:

$$\sum_{i=1}^{9n/10} f(y_i; w^T x_i) + \frac{\lambda}{2} \|w\|^2$$

10x smaller relative of weights are wrong - solutions?

just use  $w$  from one fold

multiply  $\lambda$  by  $\frac{10}{9}$  (change scale) and train on full data

Corrections to EM analysis

We showed  $\log p(x|\theta) \geq Q(\theta|\theta^*) + H(p(H|x, \theta^*))$

Part 2:  $p(x|\theta^*) = \frac{p(x, H|\theta^*)}{p(H|x, \theta^*)}$

(this was correct)

$p(x, H|\theta) = p(H|x, \theta) p(x|\theta)$

$$\sum_n p(h|x, \theta^*) \log p(x|\theta^*) = \sum_n p(h|x, \theta^*) \log p(x, h|\theta^*) - \sum_n p(h|x, \theta^*) \log p(h|x, \theta^*)$$

$$\log p(x|\theta^*) \stackrel{\sum_n p(h|x, \theta^*)}{=} Q(\theta^*|\theta^*) + H(p(H|x, \theta^*))$$

$$\Rightarrow \log p(x|\theta) - \log p(x|\theta^*) \geq Q(\theta|\theta^*) - Q(\theta^*|\theta^*) \pm H(\cdot)$$

entropy

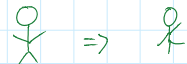
Supervised Latent Variable Models

1. Generative model with mixture for  $p(x_i|y_i)$
2. Mixtures of classifiers (mixture of experts)
3. Latent SVM

$$\min_w \frac{\lambda}{2} \|w\|^2 + \sum_{i=1}^n \max \{0, 1 - y_i \cdot w^T \phi(x_i)\}$$

$$\max_{h \in H} \{w^T \phi(x_i, h)\} \quad w^T \phi(x_i)$$

latent      standard



deformable part models

End of Scope for Midterm

## Approaches to SSL

(non-convex world)

1. EM (for generative models only)

2. Co-training

- split features into "views"

classifying webpage:  $\rightarrow$  text on webpage  
 $\rightarrow$  hyperlinks

- train a classifier on each view

$\rightarrow$  add high confidence unlabeled examples to labeled  $\rightarrow$  from either view

3. Entropy regularization

- penalize entropy of  $p(y_u | x_u, w)$

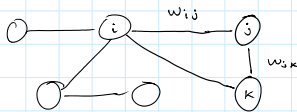
- want labels to be non-random

4. Transductive SVMs: "Just say no" (it's bad)

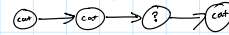
5. Graph-based SSL

- use features or relationship between  $x_i$  to define a graph

- graph has weight  $w_{ij}$ , how much we want  $y_i$  and  $y_j$  to agree



youtube example:



watching labeled videos of cats in a row, the missing labeled element is also probably a cat video

## Problems with MAP estimation

- Does MAP make the right decision?

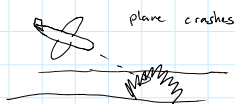
$$H = \{h_1, h_2, h_3, h_4\}$$

$h$ : hypothesis

$D$ : data

$H$ : hypothesis space

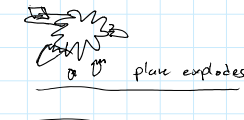
$$p(h_1 | D) = 0.25$$



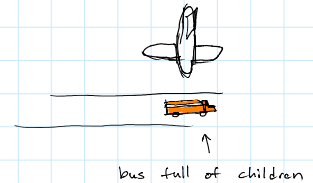
$$p(h_2 | D) = 0.3$$



$$p(h_3 | D) = 0.25$$



$$p(h_4 | D) = 0.2$$



$$p(h_2 | D) = 0.3 \quad (\text{MAP}) \quad p(\neg h_2 | D) = 0.7$$

$$p(\text{not live} | D) = p(h_1 | D) + p(h_3 | D) + p(h_4 | D) = 0.7$$

if we want to live, MAP solution doesn't exactly represent what we should do

# Learning Principles

ML:  $\hat{h} = \operatorname{argmax} p(D|h)$  train

point estimate (one  $\hat{h}$ )  $\downarrow$  predict using  $\operatorname{argmax}_{\hat{h}} p(\hat{D}, \hat{h})$  test

hidden:  $\hat{h} = \operatorname{argmax}_h \sum_z p(h, z | D)$   $\uparrow$  hidden

bagging:  $\hat{h} = \sum_{D: \epsilon \text{ bootstrap}(D)} \operatorname{argmax}_h p(h | D)$

MAP:  $\hat{h} = \operatorname{argmax} p(D|h) \propto p(D|h) p(h)$

predict using  $\operatorname{argmax}_{\hat{h}} p(\hat{D}, \hat{h})$

Bayesian: work with full posterior  $p(h|D) = \frac{p(D|h)p(h)}{p(D)} = \frac{p(D|h)p(h)}{\int p(D|h)p(h) dh}$   
 (not a "point" estimate)  
 predict by integrating over "hidden" parameters

$$\begin{aligned} p(\hat{D}|D) &= \int_H p(\hat{D}, h | D) dh \\ &= \int_H p(\hat{D} | h, D) p(h | D) dh \\ &= \int_H p(\hat{D} | h) p(h | D) dh \end{aligned}$$

## Example: Coin Flipping

Bernoulli Likelihood:

$p(X = 'H' | \theta) = \theta$   $\downarrow$  probability it comes up heads

$p(X = 'T' | \theta) = (1 - \theta)$

$p(X | \theta) = \theta^{I(X='H')} (1-\theta)^{I(X='T')}$

Beta Prior on  $\theta$ :

$\theta \sim \text{Beta}(a, b)$

$p(\theta | a, b) = \frac{1}{B(a, b)} \theta^{a-1} (1-\theta)^{b-1} \propto \theta^{a-1} (1-\theta)^{b-1}$

"Beta" function:  $B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$

uniform:  $\text{Beta}(1, 1)$

Posterior:

- assume we observe 'HHH'

$$\begin{aligned} p(\theta | \{HHH\}) &\propto \frac{p(\{HHH\} | \theta) p(\theta)}{p(\{HHH\})} \\ &= \frac{p(H|\theta) p(H|\theta) p(H|\theta) p(\theta)}{p(\{HHH\})} \\ &= \frac{\theta^3 \theta^{a-1} (1-\theta)^{b-1}}{p(\{HHH\})} \\ &= \frac{\theta^{(3+a)-1} (1-\theta)^{b-1}}{p(\{HHH\})} \end{aligned}$$

$\theta | \{HHH\} \sim \text{Beta}(3+a, b)$

Note:

we hid  $p(\{HHH\} | a, b) = B(3+a, b)$

"Marginal likelihood"

$$\text{MLE: } \theta = \frac{N_1}{N} = \frac{3}{3} = 1$$

$$\text{MAP: } \theta = \frac{\alpha - 1}{\alpha + b - 2} = \frac{3 + a - 1}{3 + a + b - 2} \xrightarrow{\text{under uniform}} \frac{3}{3} = 1$$

$$\text{Mean of posterior: } \frac{\alpha}{\alpha + b} = \frac{3 + a}{3 + a + b} \xrightarrow{\text{uniform}} \frac{4}{5} \approx 80 \text{ heads}$$

$a=3, b=3$ ; like we have  $\{HTHT\}$   
before we see data  
"expectations" on how it will behave

$$\begin{aligned} p(\hat{H} | \{HHH\}) &= \int_0^1 p(\hat{H} | \theta) p(\theta | \{HHH\}) d\theta \\ &= \int_0^1 \underbrace{\text{Ber}(\hat{H} | \theta)}_{\theta} \text{Beta}(\theta | 3+a, b) d\theta \\ &= E[\text{Beta}(\theta | 3+a, b)] \\ &= \frac{3+a}{3+a+b} \quad (\text{not generally the same as posterior mean}) \end{aligned}$$