

# Lecture 3

Note Title

1/18/2011

In the first lecture we discussed the task of labeling pixels, - either as edge or non-edge or as sky, road, vegetation, etc. We made a restrictive assumption, that we could model each pixel independently.

Suppose  $\underline{w} = \{w_i; i \in \Omega\}$  is the set of labels for each image pixel ( $\Omega = \text{set of all pixels}$ ).

$$\underline{P}(\underline{w} | I) = \prod_i P(w_i | I)$$

factorization  
conditional independence

This assumption is clearly incorrect for real images - but it is a useful approximation because it simplifies inference and learning.

Inference is simplified because to find

$$\hat{\underline{w}} = \underset{\underline{w}}{\text{ARG MAX}} P(\underline{w} | I)$$

we simply need to solve  $\hat{w}_i = \underset{w_i}{\text{ARG MAX}} P(w_i | I)$  for each  $w_i$

Note: Suppose that each pixel  $w_i$  can take  $k$  labels. Then the total number of labels for the entire image is  $k^{|\Omega|} \rightarrow$  exponentially big. (set of all possible  $\underline{w}$ )

Searching through this space of all possible labels is very hard in general.

There are exceptional cases where it is possible:

depending on the form of  $P(\underline{w} | I)$

$\rightarrow$  we may be able to use dynamic programming, steepest descent, or graph cuts. (discussed later)

If we can factorize  $P(\underline{w} | I) = \prod_i P(w_i | I)$  then we can determine  $\hat{\underline{w}}$  in at most  $|\Omega|k$  operations

Learning will also be simplified in this case. The number of training examples will be small.

(2) Now let's consider the more general case.

We can express

$$P(\underline{w} | \underline{I}, \underline{\lambda}) = \frac{1}{Z(\underline{\lambda}, \underline{I})} e^{\underline{\lambda} \cdot \underline{\phi}(\underline{w}, \underline{I})}$$

This is an exponential distribution

The  $\underline{\phi}$ 's are feature functions of  $\underline{w}$  &  $\underline{I}$  (examples later) and the  $\underline{\lambda}$ 's are parameters of the model.

$$\underline{\lambda} \cdot \underline{\phi}(\underline{w}, \underline{I}) = \sum_{\mu} \lambda_{\mu} \phi_{\mu}(\underline{w}, \underline{I})$$

Assume that we are given a set of images with the value of  $\underline{w}$  specified  $\{(\underline{w}_a, \underline{I}_a) : a \in A\}$

Learning can be formulated as finding the values of the parameters  $\underline{\lambda}$  that maximize:

$$\prod_{a \in A} P(\underline{w}_a | \underline{I}_a, \underline{\lambda})$$

Note: this is supervised learning. We assume that the potentials  $\phi_{\mu}(\underline{w}, \underline{I})$  are known, but their parameters  $\lambda_{\mu}$  are not. We will see other forms of learning in this course.

Note: this is a form of statistical regression.

Reformulate the problem in terms of minimization

$$\hat{\underline{\lambda}} = \underset{\underline{\lambda}}{\text{Arg min}} \left\{ - \sum_{a \in A} \log P(\underline{w}_a | \underline{I}_a, \underline{\lambda}) \right\}$$

Properties: (i) this is a convex optimization problem, so there is a unique solution and algorithms which are guaranteed to find it.  
(ii) the solution  $\hat{\underline{\lambda}}$  obeys.

$$\sum_{a \in A} \sum_{\underline{w}} P(\underline{w} | \underline{I}_a, \hat{\underline{\lambda}}) \underline{\phi}(\underline{w}, \underline{I}_a) = \sum_{a \in A} \underline{\phi}(\underline{w}_a, \underline{I}_a)$$

expected statistics of model = observed statistics of data.

But there is a problem. The learning algorithm may require an exponential number of operations.

(3)

Justify these two properties,

$$P(\underline{w}|I, \lambda) = \frac{e^{\lambda \cdot \underline{\phi}(\underline{w}, I)}}{Z[\lambda, I]}$$

$$-\sum_{a \in A} \log P(\underline{w}_a | I_a, \lambda) = -\sum_{a \in A} \lambda \cdot \underline{\phi}(\underline{w}_a, I_a) + \sum_a \log Z[\lambda, I_a]$$

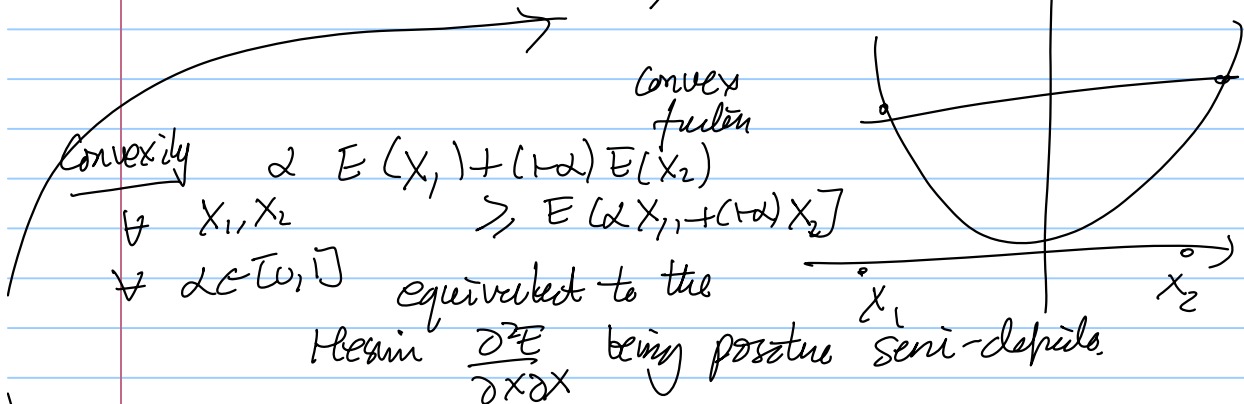
$$Z[\lambda, I_a] = \sum_{\underline{w}} e^{\lambda \cdot \underline{\phi}(\underline{w}, I_a)} \quad (\text{to ensure that } \sum_{\underline{w}} P(\underline{w}|I, \lambda) = 1)$$

The function  $\log Z[\lambda, I]$  has two important properties:

$$(a) \quad \frac{\partial}{\partial \lambda} \log Z[\lambda, I] = \sum_{\underline{w}} \underline{\phi}(\underline{w}, I) e^{\lambda \cdot \underline{\phi}(\underline{w}, I)} = \sum_{\underline{w}} \underline{\phi}(\underline{w}, I) P(\underline{w}|I, \lambda)$$

$$(b) \quad \frac{\partial^2}{\partial \lambda \partial \lambda} \log Z[\lambda, I] \text{ is positive semi-definite} \\ (\text{proof by Cauchy-Schwarz inequality})$$

Property (b) implies that  $-\sum_{a \in A} \log P(\underline{w}_a | I_a, \lambda)$  is convex and hence has, at most, one minimum



The function is bounded below by 0, so there is always one minimum.

Property (a) shows that the solution  $\hat{\lambda}$  satisfies  $\frac{\partial}{\partial \lambda} \left( -\sum_{a \in A} \log P(\underline{w}_a | I_a, \lambda) \right) = 0$

obeys

$$\sum_{a \in A} \underline{\phi}(\underline{w}_a, I_a) = \sum_{a \in A} \underline{\phi}(\underline{w}, I_a) P(\underline{w} | I_a, \hat{\lambda})$$

So far, everything looks good. But now we have to consider the computational complexity.

## (4) Learning Algorithms:

There are a variety of learning algorithms that can converge to the ~~st~~ solution  $\uparrow$ .

(1) Stochastic Descent in  $-\sum_{a \in A} \log P(\omega_a | I_a, \lambda)$

$$\lambda^{t+1} = \lambda^t - \epsilon \frac{\partial}{\partial \lambda} \left( -\sum_{a \in A} \log P(\omega_a | I_a, \lambda) \right)$$

i.e.  $\lambda^{t+1} = \lambda^t - \epsilon \sum_{a \in A} \varphi(\omega_a, I_a) + \epsilon \sum_{a \in A} \sum_w \varphi(\omega, I) P(\omega | I_a, \lambda)$

the derivative evaluates at  $\lambda^t$ .

This converges when  $\sum_{a \in A} \varphi(\omega_a, I_a) = \sum_{a \in A} \sum_w \varphi(\omega, I) P(\omega | I_a, \lambda)$   
( $\epsilon$  is time step)

(2) Generalized Iterative Scaling.

$$\lambda^{t+1} = \lambda^t - \log \sum_{a \in A} \varphi(\omega_a, I_a) + \log \sum_{a \in A} \sum_w \varphi(\omega, I) P(\omega | I_a, \lambda)$$

Discrete Iterative Algorithm

same fixed point  $\rightarrow$  note, no need of a time step  $\cdot \epsilon$ .

Important  $\rightarrow$  both methods require computing

$$\sum_{a \in A} \varphi(\omega_a, I_a) \rightarrow \text{which is practical.}$$

$$\sum_{a \in A} \sum_w \varphi(\omega, I) P(\omega | I_a, \lambda) \rightarrow \text{which is often impossible}$$

summation over  $k^{|A|}$  possible states

Two Problems: (i) normalizing to compute  $Z(\lambda) = \sum_w e^{\lambda \cdot \varphi(\omega, I)}$

(ii) computing the expectation  $\sum_w \varphi(\omega, I) P(\omega | I, \lambda)$

Again  $\rightarrow$  for some exceptional distributions these quantities can be computed (e.g. by Dynamic Programming)

But in general they cannot be.

What can we do? One possibility is Markov Chain Monte Carlo (MCMC). If we can sample efficiently from  $P(\omega | I, \lambda)$  then we compute the expectations (sampling does not require that we know  $Z(\lambda)$ .)

(5)

In standard machine learning (e.g. Stat231) we have to deal with data complexity - i.e., limited amount of data available to learn complex models.

Now, we are also faced with computational complexity - we need inference algorithms to perform tasks like:

$$\text{estimate } \hat{\underline{w}} = \underset{\underline{w}}{\text{ARG MAX}} P(\underline{w} | I, \underline{\lambda})$$

$$\text{compute } \sum_{\underline{w}} \phi(\underline{w}, I) P(\underline{w} | I, \underline{\lambda})$$

Both are possible if  $P(\underline{w} | I, \underline{\lambda})$  takes restricted forms  $\rightarrow$  e.g. distributions on graphical model without closed loops (then dynamic programming can be used).

Otherwise, we will need to use stochastic methods (MCMC) or approximate algorithms.

Note: stochastic methods are very general purpose, but they are a strategy not an algorithm. Usually a lot of problem specific knowledge is required to design an MCMC to make it computationally efficient.

How do these algorithms simplify if

$$\underline{w} = \{w_i; i \in \Lambda\}, \text{ and } P(\underline{w} | I) = \prod_i P(w_i | I)?$$

In this case, we need only learn

$$P(w_i | I, \underline{\lambda}) = \frac{1}{Z(I, \underline{\lambda}, w_i)} e^{\underline{\lambda} \cdot \phi(w_i, I)}$$

where each  $w_i$  takes only  $k$ -values

Inference is easy  $\rightarrow \hat{w}_i$  is computed by exhaustive search over  $k$ -possible values of  $w_i$ .

$\sum_{w_i} \phi(w_i, I) P(w_i | I, \underline{\lambda})$  requires  $k$  computations only.

(6) Example: suppose  $w \in \{-1, 1\}$

Define potential function  $\phi_{\mu}(w, I) = w \phi_{\mu}(I)$

$$P(w|I, \lambda) = \frac{1}{Z(I, \lambda)} e^{w \sum_{\mu} \lambda_{\mu} \phi_{\mu}(I)}$$

Learning algorithms as before. Maximize

$$\prod_{\mu} P(w^{\mu} | I^{\mu}, \lambda) \quad \text{w.r.t. } \lambda$$

The only difference is that the computations

$$\hat{w} = \underset{w}{\text{ARG MAX}} P(w|I, \lambda)$$

and  $\sum_w w \phi_{\mu}(I) P(w|I, \lambda)$  can be performed simply, because  $w$  takes only two possible values

Relationship to AdaBoost

AdaBoost seeks to learn a strong classifier

$$H(I) = \text{sign}\left(w \sum_{\mu} \lambda_{\mu} \phi_{\mu}(I)\right),$$

by selecting and combining weak classifiers

$$\{\phi_{\mu}(I) : \mu \in \Lambda\} \quad \text{with condition } \phi_{\mu}(I) \in \{\pm 1\}$$

AdaBoost can be formalized mathematically as "coordinate descent" on a function

$$F(\lambda) = \sum_{a \in A} e^{-w_a \sum_{\mu} \lambda_{\mu} \phi_{\mu}(I)}$$

(note sign change in exponential)

Initialize:  $\lambda = 0$

At iteration step  $t$

current state  $\{\lambda_{\mu}^t\}$

For each  $\mu$ , solve for  $\Delta_{\mu}$  s.t.  $\frac{\partial F}{\partial \lambda_{\mu}}(\lambda_{\mu}^t + \Delta_{\mu}) = 0$  ①

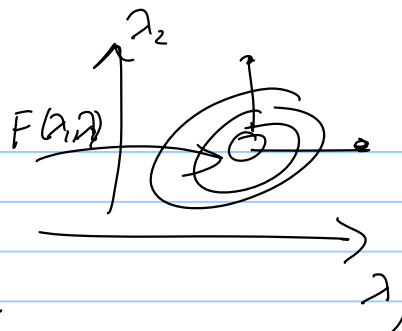
(i.e. minimize w.r.t.  $\lambda_{\mu}$  keeping the other  $\lambda$ 's fixed).

calculate  $F_{\mu}^t = F(\lambda_1^t, \dots, \lambda_{\mu}^t + \Delta_{\mu}, \dots, \lambda_n^t, \dots)$  ②

Calculate  $\hat{\mu} = \text{ARG MIN}_{\mu} F_{\mu}^t$   
set  $\lambda_{\mu}^{t+1} = \lambda_{\mu}^t$ ,  $\mu \neq \hat{\mu}$ ,  $\lambda_{\hat{\mu}}^{t+1} = \lambda_{\hat{\mu}}^t + \Delta_{\hat{\mu}}$

(7)

Geometrically,



At time  $t$ , select the direction  $\lambda_{\hat{\mu}}$  that gives maximum decrease in  $F(\lambda)$ .  
 (while not changing the other  $\lambda_{\mu}, \mu \neq \hat{\mu}$ )

Selection & Weighting:

Recall that we initialized the process so that  $\lambda_{\mu}^{t=0} = 0, \forall \mu$ .

Suppose we select  $\lambda_{\hat{\mu}}$  at time  $t+1$ , then there are two cases to consider:

(i) If  $\lambda_{\hat{\mu}}^t = 0$ , then the weak classifier  $\phi_{\hat{\mu}}(I)$  has not been used and is now selected  $\lambda_{\hat{\mu}}^{t+1} \neq 0$  and we are assigning it a weight  $\Delta_{\hat{\mu}}$ .

(ii) If  $\lambda_{\hat{\mu}}^t \neq 0$ , then we have already selected the weak classifier  $\phi_{\hat{\mu}}(I)$  and so we are simply changing its weight.  $\lambda_{\hat{\mu}}^t \rightarrow \lambda_{\hat{\mu}}^t + \Delta_{\hat{\mu}}$ .

Hence selection and weighting are performed by the same process - coordinate descent on  $F(\lambda)$

An important point to realize is that because of the form of  $F(\lambda)$  it is possible to perform all the learning calculations ①② analytically.

This reduces to the standard equation of AdaBoost.

Note: AdaBoost minimizes  $F(\lambda)$ , while conventional methods maximize  $\prod_{a \in A} P(w_a(I_a) | \lambda)$

- which is computationally harder.

But, it can be shown that both methods are asymptotically equivalent (when there is infinite data).

and that  $F(\lambda)$  &  $\prod_{a \in A} P(w_a(I_a) | \lambda)$  are related by bounds.

Note: bounds will be discussed later in the course.