

(1)

Lecture 4

Note Title

10/7/2006

Generative models $P(I|W)$, $P(W)$ How to learn $P(W)$?

For simplicity, we will discuss learning a distribution $P(W)$. Replace w by X in this lecture.

Ideal Method:

Assume a parameterized model for the distribution of form $P(X|\lambda)$ λ model parameters.

E.g. Gaussian distribution:

$$P(X|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(X-\mu)^2}{2\sigma^2}} \quad \lambda = (\mu, \sigma)$$

Assume that data is independent identically distributed (iid).

$$P(X_1, \dots, X_N | \lambda) = \prod_{i=1}^N P(X_i | \lambda) \quad (\text{product for independence})$$

Choose: $\hat{\lambda} = \underset{\lambda}{\text{ARG MAX}} P(X_1 \dots X_N | \lambda) = \text{ARG MAX} \log P(X_1 \dots X_N | \lambda)$

Hence $P(X_1, \dots, X_N | \hat{\lambda}) \geq P(X_1, \dots, X_N | \lambda)$, for all λ

(2) Example: Gaussian

$$\begin{aligned}\log P(X_1 \dots X_N | \mu, \sigma) &= \sum_{i=1}^N \log P(X_i | \mu, \sigma) \\ &= -\sum_{i=1}^N \frac{(X_i - \mu)^2}{2\sigma^2} - \sum_{i=1}^N \log \sqrt{2\pi}\sigma.\end{aligned}$$

Differentiate w.r.t. μ, σ gives

$$\frac{\partial}{\partial \mu} \log P(X_1 \dots X_N | \mu, \sigma) = \frac{1}{\sigma^2} \sum_{i=1}^N (X_i - \mu).$$

$$\frac{\partial}{\partial \sigma} \log P(X_1 \dots X_N | \mu, \sigma) = \frac{1}{\sigma^3} \sum_{i=1}^N (X_i - \mu)^2 - \frac{N}{\sigma}.$$

Maxima occur at

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N X_i$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (X_i - \hat{\mu})^2.$$

Easy to check these are maxima by computing

$$\frac{\partial^2}{\partial \mu^2}, \frac{\partial^2}{\partial \mu \partial \sigma}, \frac{\partial^2}{\partial \sigma^2}.$$

Note: similar results hold for Gaussian distributions in many variables.

Note: The Gaussian is a special case - It is often impossible to solve $\frac{\partial}{\partial \lambda} \log P(X_1 \dots X_N | \lambda) = 0$ analytically. An algorithm is required.

(3)

Exponential Distributions

$$P(\underline{x}|\underline{\lambda}) = \frac{1}{Z[\underline{\lambda}]} e^{\underline{\lambda} \cdot \underline{\phi}(\underline{x})}$$

$\underline{\lambda}$ - parameters

$\underline{\phi}(\underline{x})$ - statistics.

$$\underline{\lambda} = (\lambda_1, \lambda_2, \dots, \lambda_M)$$

$$\underline{\phi}(\underline{x}) = (\phi_1(\underline{x}), \phi_2(\underline{x}), \dots, \phi_M(\underline{x}))$$

normalization factor.

Almost every named distribution can be expressed as an exponential distribution.

For Gaussian in 1-dimension

write $\underline{\phi}(\underline{x}) = (x, x^2)$ $\underline{\lambda} = \lambda_1, \lambda_2$

$$P(x|\underline{\lambda}) = \frac{1}{Z[\underline{\lambda}]} e^{\lambda_1 x + \lambda_2 x^2}$$

compare to $\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

Translation

$$\begin{cases} \lambda_2 = -1/2\sigma^2 \\ \lambda_1 = \mu/\sigma^2 \\ Z[\underline{\lambda}] = \sqrt{2\pi}\sigma e^{\mu^2/2\sigma^2} \end{cases}$$

Similar translation into exponential distributions can be made for Poisson, Beta, Dirichlet ~ most (all) distributions you have been taught.

(4)

Learning an Exponential Distribution

You can learn them by Maximum Likelihood, which again can be interpreted in terms of minimizing the KL divergence between the empirical distribution of the data, and the model distribution.

Examples: $(x_1, x_2, \dots, x_\mu, \dots, x_N)$

$$P(\langle x_1, x_2, \dots, x_N \rangle | \lambda) = \prod_{\mu=1}^N \frac{e^{\lambda \cdot \phi(x_\mu)}}{Z(\lambda)}$$

Maximize wrt λ ||

This has a very nice form, which occurs because the exponential distribution depends on the data x only in terms of the function $\phi(x)$ - the sufficient statistics.

Note:

$$Z(\lambda) = \sum_x e^{\lambda \cdot \phi(x)} \quad \frac{\partial \log Z(\lambda)}{\partial \lambda} = \sum_x \frac{\phi(x) e^{\lambda \cdot \phi(x)}}{Z(\lambda)}$$

$$\frac{\partial \log Z(\lambda)}{\partial \lambda} = \sum_x \phi(x) P(x | \lambda)$$

(5) ML maximizes:

$$\sum_{\mu=1}^N \lambda \cdot \phi(x_{\mu}) - N \log Z[\lambda]$$

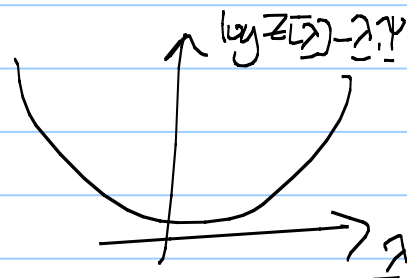
$$\frac{\partial}{\partial \lambda} \rightarrow \sum_{\mu=1}^N \phi(x_{\mu}) - N \sum_{\underline{x}} \phi(\underline{x}) P(\underline{x}|\lambda),$$

$$\sum_{\underline{x}} \phi(\underline{x}) P(\underline{x}|\lambda) = \frac{1}{N} \sum_{\mu=1}^N \phi(x_{\mu})$$

Pick the parameters λ so that the average of the statistics $\phi(\underline{x})$ w.r.t. distribution $P(\underline{x}|\lambda)$ is equal to the average of the statistics of the sampler.

Solve: $\sum_{\underline{x}} \phi(\underline{x}) P(\underline{x}|\lambda) = \frac{1}{N} \sum_{\mu=1}^N \phi(x_{\mu})$
with $\psi = \frac{1}{N} \sum_{\mu=1}^N \phi(x_{\mu})$

This is equivalent to minimizing $\log Z[\lambda] - \lambda \cdot \psi$



It can be shown that this function is convex and has a unique solution:

(Because $\frac{\partial^2}{\partial \lambda^2} (\log Z[\lambda] - \lambda \cdot \psi)$ is positive definite).

(6)

ML estimation for exponential distributions is a convex optimization function — this means that there are algorithms which are guaranteed to converge to the correct solution.

Example: Generalized Iterative Scaling (GIS)

$$\left\{ \begin{array}{l} \underline{\lambda}^{t+1} = \underline{\lambda}^t - \log \underline{\Psi}^t + \log \underline{\Psi} \\ \text{where } \underline{\Psi}^t = \sum_{\underline{x}} \phi(\underline{x}) P(\underline{x} | \underline{\lambda}^t). \\ \text{Notation: } \log \underline{\Psi} \text{ is a vector with} \\ \text{components } \log \Psi_1, \log \Psi_2, \dots, \log \Psi_m \end{array} \right.$$

But this requires computing

$$\sum_{\underline{x}} \phi(\underline{x}) P(\underline{x} | \underline{\lambda}^t)$$

which is often difficult.

(For people who took CS 202, there are often statistical / MCMC methods which can (approximately) compute this rapidly.)

(7) How does this apply to vision?

$\dots w_i \quad w_j \quad \dots$

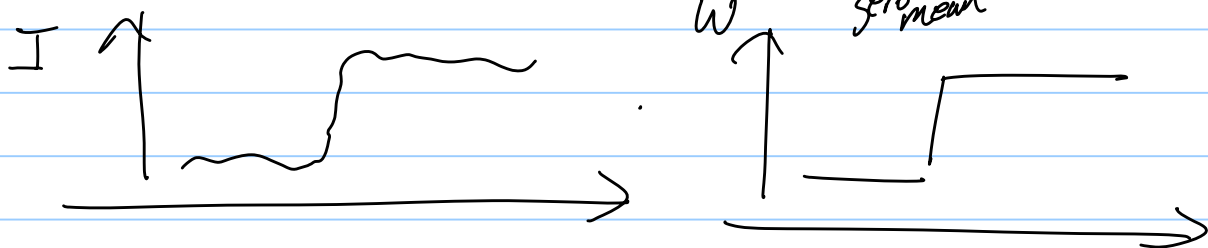
Consider the weak membrane model of images
Geman & Geman, Blake & Zisserman, Mumford & Shalek.

In probabilistic terms, this can be formulated
as $P(I|w) = \prod_i P(I_i|w_i)$
with $P(I_i|w_i) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(I_i - w_i)^2}{2\sigma^2}}$

The observed image I is a corrupted version of a
true image w . The corruption is by additive Gaussian noise

$$I_i = w_i + \epsilon_i \quad \epsilon_i \sim N(0, \sigma)$$

\uparrow noise \uparrow zero mean \uparrow variance



Need a prior $P(w)$ for the 'ideal image'.

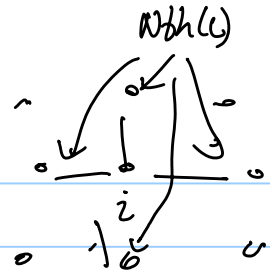
Gibbs Distribution $P(w) = \frac{1}{Z} e^{-E[w]}$

$$E[w] = \sum_i \sum_{j \in \text{nb}(i)} \psi(w_i, w_j)$$

(8)

$Nbd(i)$ is the neighborhood structure

eg. image lattice.



What function $\Psi(w_i, w_j)$?

A natural choice is $\Psi(w_i, w_j) = (w_i - w_j)^2$

penalize the square of the difference between the intensities of neighboring pixels.

Advantages:

(i) this makes it easy to learn the distribution from training data. It is a Gaussian distribution which, as we have seen, can be learnt by analytic methods (ie. no need for steepest descent or G-S).

(ii) this makes inference easy. To estimate $\hat{w} = \text{ARG MAX } P(w|I)$

reduces to minimizing.



$$E[w] = \frac{1}{2\sigma^2} \sum_i (w_i - \bar{I}_i)^2 + \sum_i \sum_{j \in Nbd(i)} \Psi(w_i, w_j)$$

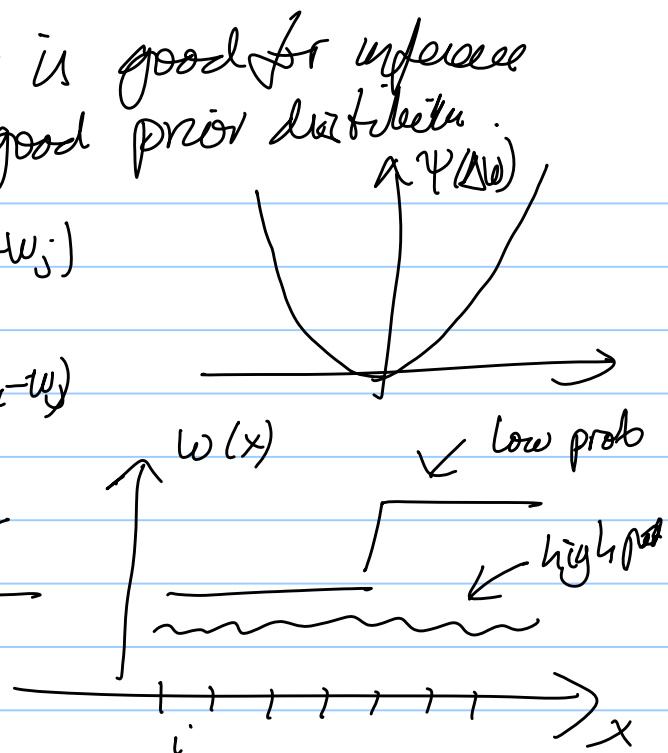
If $\Psi(w_i, w_j)$ is quadratic - $(w_i - w_j)^2$ - then $E[w]$ is quadratic, and its minimum can be found by solving linear equations $\frac{\partial E}{\partial w} = 0$

(9) But although $(w_i - w_j)^2$ is good for inference and learning it is not a good prior distribution.

It penalizes large $\Delta w = (w_i - w_j)$ too much.

It penalizes small $\Delta w = (w_i - w_j)$ too little.

It prefers images like  and dislikes images like 



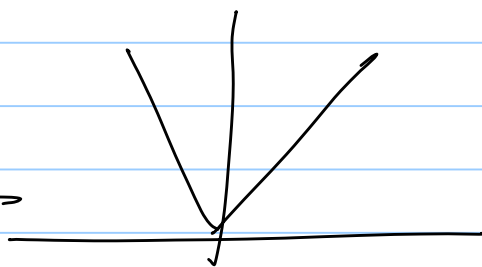
A better prior is $\Phi(w_i, w_j) = |w_i - w_j|$

this prefers images that are

very smooth 

it discourages images like 

but tolerates them better than the quadratic penalty.

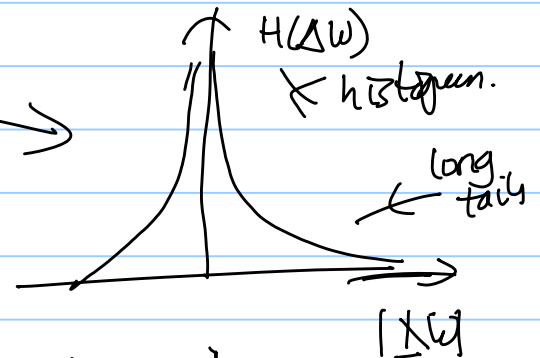


New, try to learn the prior from data.

(10)

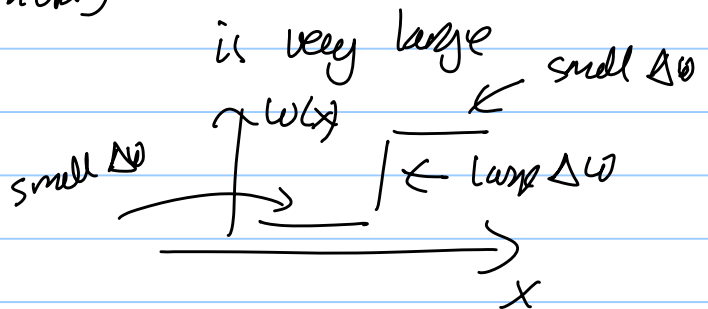
Observation, plot the histogram of $\Delta w = (w_i - w_j) \quad j \in \mathcal{N}(i)$ for image pixels.

This takes a standard form



This is inconsistent with a Gaussian distribution. \rightarrow the distribution is peaked at 0 i.e. the intensity of most pixels is similar to those of their neighbors

but Δw can also be large - i.e. at edges.



Measure the histogram

$$H(\alpha) = \sum_i \sum_{j \in \mathcal{N}(i)} \delta(\alpha, w_i - w_j)$$

$$\delta(z, w_i - w_j) = \begin{cases} 1, & \text{if } z = w_i - w_j \\ 0, & \text{otherwise.} \end{cases}$$

Exponential distribution

(max-entropy principle)
next lecture

$$-\sum_{\alpha} \lambda(\alpha) H(\alpha)$$

$$P(w) = \frac{1}{Z} e^{-\sum_{\alpha} \lambda(\alpha) H(\alpha)}$$

$$\sum_{\alpha} \lambda(\alpha) H(\alpha) = \sum_{\alpha} \lambda(\alpha) \sum_i \sum_{j \in \mathcal{N}(i)} \delta(\alpha, w_i - w_j)$$

$$= \sum_i \sum_{j \in \mathcal{N}(i)} \lambda(w_i - w_j)$$

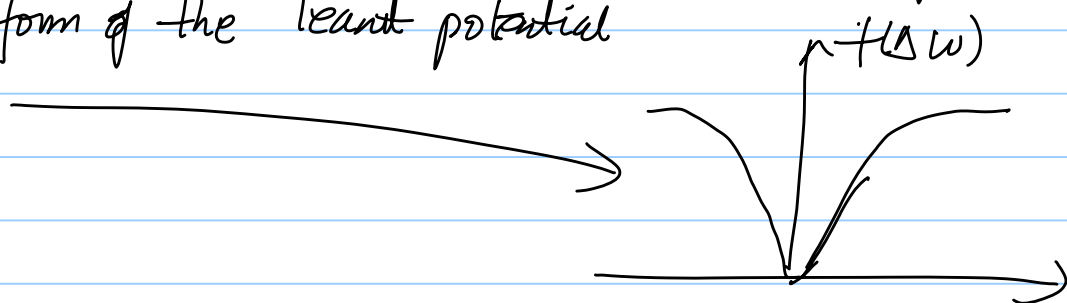
(11) Here
$$P(w) = \frac{1}{Z} e^{-\sum_i \sum_{j \in N(w_i)} \lambda (w_i - w_j)}$$

Thus, we obtain a probability distribution with the same form of the potentials - i.e. $\sum_i \sum_{j \in N(w_i)} \phi(w_i, w_j)$ -

by assuming an exponential form for the distribution and choosing the statistic $t(\lambda) = \sum_i \sum_{j \in N(w_i)} \delta(\lambda, w_i - w_j)$.

(Zhu & Hampel (1997))

The form of the learnt potential is



This was suggested earlier (Geman & Geman) by specifying additional "line process" variables.

$$P(w) \rightarrow P(w, L) = \frac{1}{Z} e^{-E(w, L)}$$

$$E(w, L) = \sum_i \sum_{j \in N(w_i)} (w_i - w_j)^2 (1 - L_{ij}) + \kappa \sum_i \sum_{j \in N(w_i)} L_{ij}$$

$L_{ij} \in \{0, 1\}$, is a binary-valued

variable. $L_{ij} = 1$ 'cuts' the smoothness between w_i & w_j .

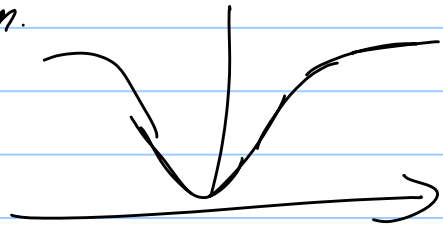
Lower Energy (higher probability) if $L_{ij} = 1$, when $|w_i - w_j| > \kappa^{\frac{1}{2}}$
 $L_{ij} = 0$, when $|w_i - w_j| < \kappa^{\frac{1}{2}}$

(12)

$$P(w, L) = \frac{1}{Z} e^{-E(w, L)}$$

$$P(w) = \sum_L P(w, L) = \frac{1}{Z} e^{-\sum_i \sum_{j \in \text{nb}(i)} \phi(w_i, w_j)}$$

It can be shown that $\phi(w_i, w_j)$ has form
- similar to the least form.



Note: more accurate priors can be obtained by
considering higher order filters $\phi(w_i, w_j, \dots, w_k)$

It can be shown that the statistics of any

filter.

$$\sum_i A_i w_i, \text{ such that } \sum_i A_i = 0$$

has very similar form on each image (M. Green, Maths, UCL)

But this leads to complicated probability distributions.

Hard to do inferences on them.