

① Derivation of IPF

Parameter learning in Undirected Models
D: dataset of m iid data points

$$\ell(D; \theta) = \log \prod_{j=1}^m \frac{1}{Z} \prod_{i=1} \psi_i(C_i = c_i^{(j)})$$

$$= \sum_{j=1}^m \underbrace{\sum_i \log \psi_i(C_i = c_i^{(j)})}_{\text{swap}} - m \log Z$$

$$= \sum_{i=1} \underbrace{\sum_{j=1}^m \log \psi_i(C_i = c_i^{(j)})}_{\text{count}} - m \log Z$$

$$= \sum_{i=1} \sum_{c_i} \underbrace{\text{count}(C_i = c_i)}_{\hat{P}} \log \psi_i(C_i = c_i) - m \log Z$$

$$= m \sum_{i=1} \underbrace{\sum_{c_i} \tilde{P}(C_i = c_i) \log \psi_i(C_i = c_i)}_{\text{Does this look like something familiar?}}$$

Does this look like something familiar?

Hint: Expectation

want to learn

Parameter: $\Psi_k(c_k)$

what does this mean for HW 5?

$$\theta = \begin{bmatrix} \theta_1 & \theta_2 \\ \theta_3 & \theta_4 \end{bmatrix}$$

Take derivative:

$$\frac{\partial \ell(D; \theta)}{\partial \Psi_k(c_k)} = \frac{\partial}{\partial \Psi_k(c_k)} \left[m \sum_{i=1}^m \sum_{c_i} \hat{P}(c_i) \log \Psi_i(c_i) \right] - \frac{\partial}{\partial \Psi_k(c_i)} [m \log Z]$$

First term $\frac{\partial}{\partial \Psi_k(c_k)} \left[m \sum_{i=1}^m \sum_{c_i} \hat{P}(c_i) \log \Psi_i(c_i) \right]$

Parameter independence assumption

only when they agree on i
ie $k=i$
does a term survive

$$= m \hat{P}(c_i) \frac{\partial}{\partial \Psi_k(c_k)} \log \Psi_k(c_k)$$

$$= m \frac{\tilde{P}(c_i)}{\Psi_k(c_k)}$$

Now

$$\frac{\partial}{\partial \Psi_k(c_k)} \log Z$$

$$= \frac{\partial}{\partial \Psi_k(c_k)} \log \sum_{\mathcal{X}} \prod_i \Psi_i(x_i)$$

$$= \frac{1}{Z} \frac{\partial}{\partial \Psi_k(c_k)} \left[\sum_{\mathcal{X}} \prod_i \Psi_i(x_i) \right]$$

$$= \frac{1}{Z} \left[\sum_{\mathcal{X}} \left\{ \prod_{i \neq k} \Psi_i(x_i) \right\} \frac{\partial \Psi_k(c_k)}{\partial \Psi_k(c_k)} \right]$$

$\underbrace{\hspace{10em}}_{= I(\mathcal{X}_k = c_k)}$

$$= \frac{1}{\Psi_k(c_k)} \left[\frac{1}{Z} \sum_{\mathcal{X}} \prod_i \Psi_i(x_i) I(\mathcal{X}_k = c_k) \right]$$

$$= \frac{1}{\Psi_k(c_k)} \left[\sum_{\mathcal{X}} \underbrace{\frac{1}{Z} \prod_i \Psi_i(x_i)}_{P(\mathcal{X})} I(\mathcal{X}_k = c_k) \right]$$

$$= \frac{E[I(\mathcal{X}_k = c_k)]}{\Psi_k(c_k)} = \frac{P(\mathcal{X}_k = c_k)}{\Psi_k(c_k)}$$

② Parameter Learning for Log-Linear Models

$$P(x) = \frac{1}{Z} \exp \left\{ \sum_i w_i \phi_i(x_i) \right\}$$

← Again D : dataset of m iid samples

$$\text{LCD}(\theta) = \log \prod_{j=1}^N \frac{1}{Z} \exp \left\{ \sum_i w_i^T \phi_i(x_i) \right\}$$

$$= \sum_j \sum_i w_i^T \phi_i(x_i) - m \log Z$$

↻
swap

$$= \sum_i \sum_j w_i^T \phi_i(x_i) - m \log Z$$

↵
count

$$= m \sum_i \sum_{c_i} \tilde{P}(x_i=c_i) w_i^T \phi_i(x_i=c_i) - m \log Z$$

Want to learn: \vec{w}_k

Take derivative:

$$\frac{\partial l(D; \theta)}{\partial \vec{w}_k} = \frac{\partial}{\partial \vec{w}_k} \left[m \sum_i \sum_{c_i} \hat{P}_k(x_i = c_i) \omega_i^T \phi_i(x_i = c_i) \right]$$

non-zero only when $k=i$

$$- \frac{\partial m \log Z}{\partial \vec{w}_k}$$

$$= m \sum_{c_i} \hat{P}(x_i = c_i) \phi_i(x_i = c_i) - \frac{\partial \log Z}{\partial \omega_i}$$

$$= m E_P[\phi_i(x_i)] - m \frac{\partial \log Z}{\partial \omega_i}$$

←————→

$$\frac{\partial \log Z}{\partial \vec{w}_k} = \frac{\partial}{\partial \omega_k} \log \sum_x \exp \left\{ \sum_i \omega_i^T \phi_i(x_i) \right\}$$
$$= \frac{1}{Z} \left[\frac{\partial}{\partial \vec{w}_k} \sum_x \exp \left\{ \sum_i \omega_i^T \phi_i(x_i) \right\} \right]$$
$$= \frac{1}{Z} \sum_x \exp \left\{ \sum_i \omega_i^T \phi_i(x_i) \right\} \frac{\partial}{\partial \omega_i} \sum_i \omega_i^T \phi_i(x_i)$$
$$= \sum_x P_k(x) \phi_i(x_i)$$
$$= E_P[\phi_i(x_i)]$$

Thus gradient

$$= m \left\{ \underbrace{E_P [\Phi_i(x_i)]}_{\substack{\text{what data} \\ \text{believes / knows} \\ \text{about features} \\ \text{distribution}}} - \underbrace{E_P [\Phi_i(x_i)]}_{\substack{\text{what our} \\ \text{model} \\ \text{believes} \\ \text{about features} \\ \text{distribution}}} \right\}$$

B Algorithm

1) Initialize weights

2) Compute gradient

3) take small step / whatever heuristic you like

4) Repeat

Again notice: Learning requires inference