# Stat 521A
# Lecture 25

# Outline

- MAP param estimation for UGMs (20.1-20.4)
- Learning using approximate inference (20.5)
- Alternative objectives (20.6)

- ## Log-linear model

$$P(X_1, \ldots, X_n : \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp \left\{ \sum_{i=1}^{k} \theta_i f_i[D_i] \right\}.$$

$$\ln Z(\boldsymbol{\theta}) = \ln \sum_{\xi} \exp \left\{ \sum_{i} \theta_i f_i[\xi] \right\}.$$
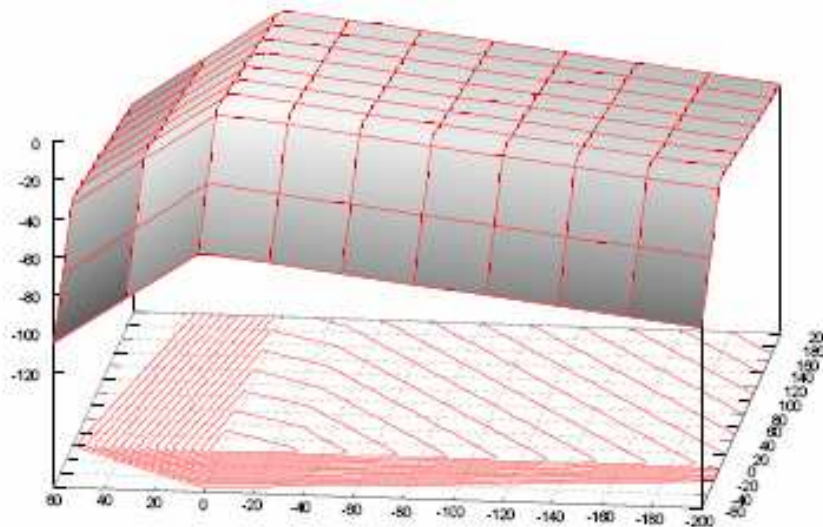
Concave



**Figure 20.1** Log-likelihood surface for the Markov network $A$—$B$—$C$, as a function of $\ln \phi_1[a^1, b^1]$ ($x$-axis) and $\ln \phi_2[b^1, c^1]$ ($y$-axis); all other parameters in both potentials are set to 1. The data set $\mathcal{D}$ has $M = 100$ instances, for which $M[a^1, b^1] = 90$ and $M[b^1, c^1] = 15$. (The other sufficient statistics are irrelevant, as all of the other log-parameters are 0.)

3

# LogZ: first deriv

**Proposition 20.2.3:** *Let $\Phi$ be a set of features. Then,*

$$\frac{\partial}{\partial \theta_i} \ln Z(\theta) = E_{\boldsymbol{\theta}}[f_i]$$

$$\frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln Z(\theta) = Cov_{\boldsymbol{\theta}}[f_i; f_j],$$

*where $E_{\boldsymbol{\theta}}[f_i]$ is a shorthand for $E_{P(\mathcal{X}:\boldsymbol{\theta})}[f_i]$.*

$$\frac{\partial}{\partial \theta_i} \ln Z(\theta) = \frac{1}{Z(\theta)} \sum_{\xi} \frac{\partial}{\partial \theta_i} \exp \left\{ \sum_j \theta_j f_j[\xi] \right\}$$

$$= \frac{1}{Z(\theta)} \sum_{\xi} f_i[\xi] \exp \left\{ \sum_j \theta_j f_j[\xi] \right\}$$

$$= \mathbb{E}_{\boldsymbol{\theta}}[f_i].$$

# logZ: second deriv

$$
\begin{aligned}
\frac{\partial^2}{\partial\theta_j\partial\theta_i}\ln Z(\theta) &= \frac{\partial}{\partial\theta_j}\left[\frac{1}{Z(\theta)}\sum_\xi f_i[\xi]\exp\left\{\sum_k \theta_k f_k[\xi]\right\}\right] \\
&= -\frac{1}{Z(\theta)^2}\left(\frac{\partial}{\partial\theta_j}Z(\theta)\right)\sum_\xi f_i[\xi]\exp\left\{\sum_k \theta_k f_k[\xi]\right\} \\
&\quad +\frac{1}{Z(\theta)}\sum_\xi f_i[\xi]f_j[\xi]\exp\left\{\sum_k \theta_k f_k[\xi]\right\} \\
&= -\frac{1}{Z(\theta)^2}Z(\theta)E_\theta[f_i]\sum_\xi f_i[\xi]\tilde{P}(\xi:\theta) \\
&\quad +\frac{1}{Z(\theta)}\sum_\xi f_i[\xi]f_j[\xi]\tilde{P}(\xi:\theta) \\
&= E_\theta[f_i]\sum_\xi f_i[\xi]P(\xi:\theta) \\
&\quad +\sum_\xi f_i[\xi]f_j[\xi]P(\xi:\theta) \\
&= E_\theta[f_i f_j] - E_\theta[f_i]E_\theta[f_j] \\
&= Cov_\theta[f_i; f_j].
\end{aligned}
$$

5

# Finding the MLE

At optimum, model moments = empirical moments

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta : \mathcal{D}) = \mathbb{E}_{\mathcal{D}}[f_i[\mathcal{X}]] - E_{\theta}[f_i]. \tag{20.4}$$

This analysis provides us with a precise characterization of the maximum likelihood parameters $\hat{\theta}$:

**Theorem 20.3.1:** *Let $\Phi$ be a set of features. Then, $\theta$ is a maximal likelihood parameter assignment if and only if $E_{\mathcal{D}}[f_i[\mathcal{X}]] = E_{\hat{\theta}}[f_i]$ for all $i$.*

Must perform inference once per gradient

Just do gradient based optimization, eg stochastic gradient descent.
Expensive to compute Hessian explicitly, so use Quasi-Newton.

$$\frac{\partial}{\partial \theta_i \partial \theta_j} \ell(\theta : \mathcal{D}) = -M\mathbb{C}\text{ov}_{\theta}[f_i; f_j].$$

# CRFs

- Conditional density models

$$\ell_{\boldsymbol{Y}|\boldsymbol{X}}(\boldsymbol{\theta} : \mathcal{D}) = \ln P(y[1,\ldots,M] \mid x[1,\ldots,M],\boldsymbol{\theta}) = \sum_{m=1}^{M} \ln P(y[m] \mid x[m],\boldsymbol{\theta}).$$

$$\frac{\partial}{\partial \theta_i} \ell_{\boldsymbol{Y}|\boldsymbol{X}}(\boldsymbol{\theta} : \mathcal{D}) = \sum_{m=1}^{M} \left[ f_i[y[m], x[m]] - E_{\boldsymbol{\theta}}[f_i \mid x[m]] \right].$$

Must perform inference M times per gradient

# MRFs with hidden variables

- Must perform inference M times per gradient

$$\frac{1}{M} \ln P(\mathcal{D} \mid \boldsymbol{\theta}) = \frac{1}{M} \ln \left( \sum_{m=1}^{M} \sum_{\boldsymbol{h}[m]} P(o[m], h[m] \mid \boldsymbol{\theta}) \right)$$

$$= \frac{1}{M} \ln \left( \sum_{m=1}^{M} \sum_{\boldsymbol{h}[m]} \tilde{P}(o[m], h[m] \mid \boldsymbol{\theta}) \right) - \ln Z.$$

$$\frac{\partial}{\partial \theta_i} \ln \sum_{\boldsymbol{h}[m]} \tilde{P}(o[m], h[m] \mid \boldsymbol{\theta}) = E_{\boldsymbol{h}[m] \sim P(\mathcal{H}[m] \mid o[m], \boldsymbol{\theta})} [f_i],$$

**Proposition 20.3.3:** *For a data set* $\mathcal{D}$

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\boldsymbol{\theta} : \mathcal{D}) = \frac{1}{M} \left[ \sum_{m=1}^{M} \mathbb{E}_{\boldsymbol{h}[m] \sim P(\mathcal{H}[m] \mid o[m], \boldsymbol{\theta})} [f_i] \right] - E_{\boldsymbol{\theta}}[f_i].$$

clamped            unclamped

# CRFs with hidden variables

- Training is similar to MRFs with hidden variables, except expectations condition on x_n, so need to be redone for each case
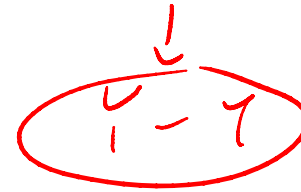
# Summary

MRF

$$\nabla = \sum_i f(x_i) - ME_X[f(X)]$$

CRF

$$\nabla = \sum_i f(x_i, y_i) - \sum_i E_Y[f(x_i, Y)]$$

MRF + H

$$\nabla = \sum_i E_H f(H_i, x_i) - ME_{H,X}[f(H, X)]$$

CRF + H

$$\nabla = \sum_i E_H f(x_i, y_i, H) - \sum_i E_{H,Y}[f(x_i, Y, H)]$$

10

# ML and MaxEnt

- MLE in the expfam is equivalent to MaxEnt subject to moment constraints

Maximum-Entropy

Find $Q(\mathcal{X})$
that maximize $H_Q(\mathcal{X})$

subject to

$$E_Q[f_i] = E_\mathcal{D}[f_i] \quad i = 1, \ldots, k$$

**Theorem 20.3.4:** *The distribution $Q^*$ is the maximum entropy distribution satisfying Eq. (20.10) if and only if $Q^* = P_{\hat{\theta}}$, where*

$$P_{\hat{\theta}}(\mathcal{X}) = \frac{1}{Z(\hat{\theta})} \exp\left\{\sum_i \hat{\theta}_i f_i[\mathcal{X}]\right\}$$

*and $\hat{\theta}$ is the maximum likelihood parameterization relative to $\mathcal{D}$.*

# Proof

PROOF For notational simplicity, let $P = P_{\hat{\theta}}$. From Theorem 20.3.1, it follows that $E_P[f_i] = E_{\mathcal{D}}[f_i[\mathcal{X}]]$ for $i = 1, \ldots, k$, and hence that $P$ satisfies the constraints of Eq. (20.10). Therefore, to prove that $P = Q^*$, we need only show that $H_P(\mathcal{X}) \geq H_Q(\mathcal{X})$ for all other distributions $Q$ that satisfy these constraints. Consider any such distribution $Q$.

From Theorem 8.4.1, it follows that:

$$H_P(\mathcal{X}) = -\sum_i \hat{\theta}_i E_P[f_i] + \ln Z(\boldsymbol{\theta}). \tag{20.11}$$

Thus,

$$
\begin{aligned}
H_P(\mathcal{X}) - H_Q(\mathcal{X}) &= -\left[\sum_i \theta_i E_P[f_i[\mathcal{X}]]\right] + \ln Z_P - E_Q[-\ln Q(\mathcal{X})] \\
(i) &= -\left[\sum_i \theta_i E_Q[f_i[\mathcal{X}]]\right] + \ln Z_P + E_Q[\ln Q(\mathcal{X})] \\
&= E_Q[-\ln P(\mathcal{X})] + E_Q[\ln Q(\mathcal{X})] \\
&= D(Q\|P) \geq 0,
\end{aligned}
$$

where (i) follows from the fact that both $P_{\hat{\theta}}$ and $Q$ satisfy the constraints, so that $E_{P_{\hat{\theta}}}[f_i] = E_Q[f_i]$ for all $i$.
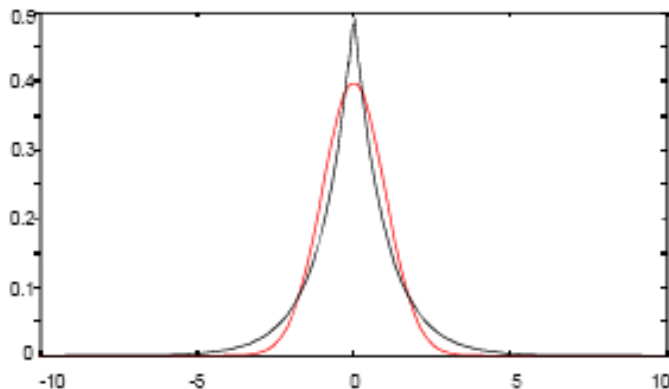
We conclude that $H_{P_{\hat{\theta}}}(\mathcal{X}) \geq H_Q(\mathcal{X})$ with equality if and only if $P_{\hat{\theta}} = Q$. Thus, the maximum entropy distribution $Q^*$ is necessarily equal to $P_{\hat{\theta}}$, proving the result. ∎

# MAP estimation

- Convex prior + convex likelihood makes objective strictly convex (unique soln)
- Also helps prevent overfitting
- L2 and L1

$$P(\theta \mid \sigma^2) = \prod_{i=1}^{k} \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{\theta_i^2}{2\sigma^2}\right\},$$

$$P_{Laplacian}(\theta \mid \beta) = \frac{1}{2\beta} \exp\left\{-\frac{|\theta|}{\beta}\right\}.$$

$$
\begin{aligned}
\ln \frac{P(\xi)}{P(\xi')} &= \sum_{i=1}^{k} \theta_i f_i[\xi] - \sum_{i=1}^{k} \theta_i f_i[\xi'] \\
&= \sum_{i=1}^{k} \theta_i (f_i[\xi] - f_i[\xi']).
\end{aligned}
$$

# Learning with approximate inference

- Recall that the gradient requires model expectation over the features

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta : \mathcal{D}) = \mathbb{E}_{\mathcal{D}}[f_i[\mathcal{X}]] - \mathbb{E}_{\theta}[f_i].$$ (20.4)

- We can use approximate inference to approximate the expectation, but approximate gradients can cause learning to diverge

# Pseudo moment matching

- At the optimum, the pseudo marginals must satisfy

$$E_{\beta_i[C_i]}[f_{C_i}] = E_{\mathcal{D}}[f_i[C_i]].$$

- Suppose we use tabular features. Then

$$\beta_i[c_i^j] = \hat{P}(c_i^j).$$

- Hence we don't need to run inference. There are multiple potentials that can generate these beliefs. We can uniquely recover one set using (for any ordering i<j)

$$\phi_i \leftarrow \frac{\beta_i}{\mu_{i,j}}.$$

# Unified inference and learning

- Pseudo moment matching only works for unconditional, tabular potentials with no tying and no regularizer

- To combine BP with param optimization, we can optimize

Approx-Maximum-Entropy

Find $Q$

that maximize $\sum_{C_i \in \mathcal{U}} H_{\beta_i}(C_i) - \sum_{(C_i, C_j) \in \mathcal{U}} H_{\mu_{i,j}}(S_{i,j})$

subject to $\quad E_{\beta_i}[f_i] = E_{\mathcal{D}}[f_i] \quad i = 1, \ldots, k$

$Q \in Local[\mathcal{U}]$

The model parameters theta are the Lagrange multiplers for E[f]
And the messages are the Lagrange multipliers for the local consistency

Handwritten notes:

$c_1$ (diagram with circle containing $A \to B$, $V$, $C$, labeled $c_2$, $c_3$)

$$A \quad B \, C$$
$$0 \quad 0 \; 0$$
$$0 \quad 1 \; 0$$
$$1 \quad 0 \; 0$$

Tied features

$f_{00}(x,y) = 1$ iff $x = y = 0$

$f_{11}(x,y) = 1$ iff $x = y = 1$

$E_0 f_{00} = I(A_1 = B_1) + I(A_1 = C_1) + I(B_1 = C_1)$
$+ \cdots = 5$

Find $Q$
that maximize $H_{\beta_1}(A, B) + H_{\beta_2}(B, C) + H_{\beta_3}(A, C) - H_{\mu_{1,2}}(B) - H_{\mu_{2,3}}(C) - H_{\mu_{2,3}}(A)$

subject to

$$\sum_i E_{\beta_i}[f_{00}] = 3 \quad (20.17)$$

$$\sum_i E_{\beta_i}[f_{11}] = 0 \quad (20.18)$$

$$\sum_a [\beta_1[a, b]] - \sum_c [\beta_2[b, c]] = 0 \quad (20.19)$$

$$\sum_b [\beta_2[b, c]] - \sum_a [\beta_3[a, c]] = 0 \quad (20.20)$$

$$\sum_c [\beta_3[a, c]] - \sum_b [\beta_1[a, b]] = 0 \quad (20.21)$$

$$\sum_{c_i} \beta_i[c_i] = 1 \quad i = 1, 2, 3 \quad (20.22)$$

$$\beta_i \geq 0 \quad i = 1, 2, 3 \quad (20.23)$$

18

# Double loop algorithm

- Inner loop optimizes $\delta_{ij}$ by iterating the fixed point eqns

- Outer loop optimizes $\theta$ eg using gradient descent

# Approximating Z

- Loglik

$$\ell(\theta : \xi) = \ln \tilde{P}(\xi \mid \theta) - \ln Z(\theta)$$

$$\ln \tilde{P}(\xi \mid \theta) - \ln \left( \sum_{\xi'} \tilde{P}(\xi' \mid \theta) \right).$$

- We can approximate the sum in different ways

# Pseudolikeliood

- Define

$$P(\xi) = \prod_{j=1}^{n} P(x_j \mid x_1, \ldots, x_{j-1}) \qquad P(\xi) \approx \prod_j P(x_j \mid x_1, \ldots, x_{j-1}, x_{j+1}, \ldots, x_n)$$

$$\ell_{\text{pseudo}}(\theta : \mathcal{D}) = \frac{1}{M} \sum_m \sum_j \ln P(x_j[m] \mid x_{-j}[m], \theta)$$

$$P(x_j \mid x_{-j}) = \frac{P(x_j, x_{-j})}{P(x_{-j})} = \frac{\tilde{P}(x_j, x_{-j})}{\tilde{P}(x_{-j})}$$

$$= \frac{\tilde{P}(x_j, x_{-j})}{\sum_{x'_i} \tilde{P}(x'_j, x_{-j})}.$$

# Gradient of PL

$\ln P(x_j \mid x_{-j}) =$                 Convex

$$\left( \sum_{i\,:\,Scope[f_i] \ni X_j} \theta_i f_i[x_j, u_j] \right) - \ln \left( \sum_{x'_j} \exp \left\{ \sum_{i\,:\,Scope[f_i] \ni X_j} \theta_i f_i[x'_j, u_j] \right\} \right).$$

$$\frac{\partial}{\partial \theta_i} \ln P(x_j \mid x_{-j}) = f_i[x_j, x_{-j}] - \mathbb{E}_{x'_j \sim P_{\boldsymbol{\theta}}(X_j \mid \boldsymbol{x}_{-j})} \left[ f_i[x'_j, x_{-j}] \right].$$

**Proposition 20.6.1:**

$$\frac{\partial}{\partial \theta_i} \ell_{\text{pseudo}}(\boldsymbol{\theta} : \mathcal{D}) = \sum_{j:X_j \in Scope[f_i]} \left( \frac{1}{M} \sum_m f_i[\xi[m]] - E_{x'_j \sim P_{\boldsymbol{\theta}}(X_j \mid \boldsymbol{x}_{-j}[m])} \left[ f_i[x'_j, x_{-j}[m]] \right] \right).$$

$$(20.31)$$

# Consistency of PL

- Thm 20.6.2 (Besag). If data is generated from our model with params θ*, then as M->inf, argmax PL(θ) -> θ*.

- Pf. The empirical approaches P(θ*). Hence

$$\frac{1}{M} \sum_m f_i[\xi[m]] \longrightarrow \mathbb{E}_{\xi \sim P_{\theta^*}(\mathcal{X})}[f_i[\xi]].$$

- And

$$\frac{1}{M} \sum_m \mathbb{E}_{x'_j \sim P_{\theta^*}(X_j | \boldsymbol{x}_{-j}[m])}[f_i[x'_j, x_{-j}[m]]] \quad = \quad \sum_{\boldsymbol{x}_{-j}} P_D(x_{-j}) \sum_{x'_j} P_{\theta^*}(x'_j \mid x_{-j}) f_i[x'_j, x_{-j}]$$

$$\longrightarrow \quad \sum_{\boldsymbol{x}_j} P_{\theta^*}(x_{-j}) \sum_{x'_j} P_{\theta^*}(x'_j \mid x_{-j}) f_i[x'_j, x_{-j}]$$

$$= \quad \mathbb{E}_{\xi \sim P_{\theta^*}}[f_i[\xi]].$$

- Hence gradient of PL is zero at θ*.

- Ex 20.6.3 (cf Hinton's greek vase)



Assume X1, X2 are strongly correlated (eg mirror images),
And X1,Y and X2,Y are less strongly correlated.
PL will learn that X1 can be predicted from X2, and will ignore Y.
At test time, if we observe Y and want to predict X1, we are hosed.

# Sample based learning

- Recall loglik is

$$\frac{1}{M}\ell(\theta : \mathcal{D}) = \sum_i \theta_i E_{\mathcal{D}}[f_i[d_i]] - \ln Z(\theta),$$

$$
\begin{aligned}
Z(\theta) &= \sum_\xi \exp\left\{\sum_i \theta_i f_i[\xi]\right\} \\
&= \sum_\xi \frac{Q(\xi)}{Q(\xi)} \exp\left\{\sum_i \theta_i f_i[\xi]\right\} \\
&= E_Q\left[\frac{1}{Q(\mathcal{X})} \exp\left\{\sum_i \theta_i f_i[\mathcal{X}]\right\}\right].
\end{aligned}
$$

Sample K x's given θ
Compute ln Z(θ)
Update θ
Repeat

$$
\begin{aligned}
Z(\theta) &= E_{P_{\theta^0}}\left[\frac{Z(\theta^0)\exp\left\{\sum_i \theta_i f_i[\mathcal{X}]\right\}}{\exp\left\{\sum_i \theta_i^0 f_i[\mathcal{X}]\right\}}\right] \\
&= Z(\theta^0) E_{P_{\theta^0}}\left[\exp\left\{\sum_i (\theta_i - \theta_i^0) f_i(\mathcal{X})\right\}\right].
\end{aligned}
$$

$$\ln Z(\theta) \approx \ln\left(\frac{1}{K}\sum_{k=1}^K \exp\left\{\sum_i (\theta_i - \theta_i^0) f_i(\xi^k)\right\}\right) + \ln Z(\theta^0).$$

# Contrastive divergence

- Might need to sample many x's to accurately approximate Z, but this is slow
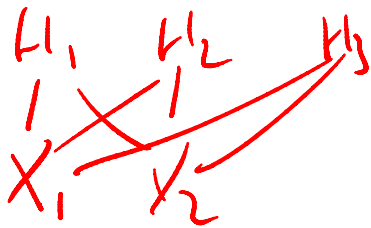- So define a set D- of randomly perturbed neighbors of D, and use

$$\ell_{\text{CD}}(\theta : \mathcal{D} \| \mathcal{D}^-) = \left[ \mathbb{E}_{\xi \sim \hat{P}_{\mathcal{D}}}[\ln P_\theta(\xi)] - \mathbb{E}_{\xi \sim \hat{P}_{\mathcal{D}-}}[\ln P_\theta(\xi)] \right],$$

$$\frac{\partial}{\partial \theta_i} \ell_{\text{CD}}(\theta : \mathcal{D} \| \mathcal{D}^-) = \mathbb{E}_{\hat{P}_{\mathcal{D}}}[f_i[\mathcal{X}]] - \mathbb{E}_{\hat{P}_{\mathcal{D}-}}[f_i].$$

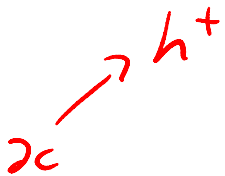- Often $x_i$- is generated by applying 1 step of Gibbs sampling to $x_i$

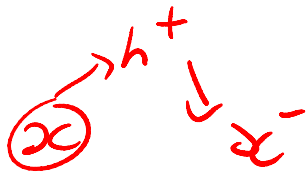- RBMs have 1 layer of hidden variables, so we need an additional expectation

$$\nabla_i = E_{\mathbf{x}\sim D}E_{\mathbf{h}}f_i(\mathbf{x},\mathbf{h}) - E_{\mathbf{x}\sim D^-}E_{\mathbf{h}}f_i(\mathbf{x},\mathbf{h})$$

$$\approx \frac{1}{N}\sum_n E_{\mathbf{h}_n}f_i(\mathbf{x}_n,\mathbf{h}_n) - E_{\mathbf{h}_n}f_i(\mathbf{x}_n^-,\mathbf{h}_n)$$

$$\approx \frac{1}{N}\sum_n f_i(\mathbf{x}_n,\mathbf{h}_n^+) - f_i(\mathbf{x}_n^-,\mathbf{h}_n^-)$$
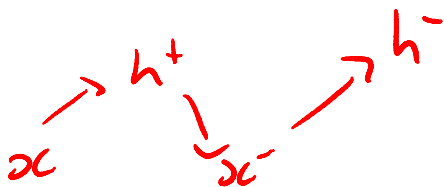
Stop learning when your dreams match reality

$$\mathbf{h}_n^+ \sim p(\mathbf{h}|\mathbf{x}_n,\boldsymbol{\theta}) \qquad \text{Interpretation of data}$$

$$\mathbf{x}_n^- \sim p(\mathbf{x}|\mathbf{h}_n^+,\boldsymbol{\theta}) \qquad \text{Reconstruction/Fantasy data}$$

$$\mathbf{h}_n^- \sim p(\mathbf{h}|\mathbf{x}_n^-,\boldsymbol{\theta}) \qquad \text{Interpretation of your fantasies}$$

28

# MAP approximation (perceptron training)

- Let us approximate Z (sum over all X) by the MAP estimate. Objective becomes

$$\frac{1}{M}\ell(\theta : \mathcal{D}) - \ln P(\xi^{\mathrm{MAP}}(\theta) \mid \theta), \qquad \frac{1}{M}\sum_{m=1}^{M} \ln \tilde{P}(\xi[m] \mid \theta) - \ln \tilde{P}(\xi^{\mathrm{MAP}}(\theta) \mid \theta).$$

- For a single data term

$$\begin{aligned}
\ln P(\xi[m] \mid \theta) &- \ln P(\xi^{\mathrm{MAP}}(\theta) \mid \theta) \\
&= [\ln \tilde{P}(\xi[m] \mid \theta) - \ln Z(\theta)] - [\ln \tilde{P}(\xi^{\mathrm{MAP}}(\theta) \mid \theta) - \ln Z(\theta)] \\
&= \ln \tilde{P}(\xi[m] \mid \theta) - \ln \tilde{P}(\xi^{\mathrm{MAP}}(\theta) \mid \theta) \\
&= \sum_i \theta_i [f_i[\xi[m]] - f_i[\xi^{\mathrm{MAP}}(\theta)]].
\end{aligned}$$

- Hence gradient is

$$E_{\mathcal{D}}[f_i[\mathcal{X}]] - f_i[\xi^{\mathrm{MAP}}(\theta)].$$

# Problem with MAP approximation

- The objective is always negative or 0 since

$$\frac{1}{M}\ell(\boldsymbol{\theta} : \mathcal{D}) - \ln P(\xi^{\text{MAP}}(\boldsymbol{\theta}) \mid \boldsymbol{\theta}),$$

$$\ln P(\xi[m] \mid \boldsymbol{\theta}) \leq \ln P(\xi^{\text{MAP}}(\boldsymbol{\theta}) \mid \boldsymbol{\theta}),$$

- We can always achieve the maximum of 0 by setting \theta=0

$$
\begin{aligned}
\ln P(\xi[m] \mid \boldsymbol{\theta}) &- \ln P(\xi^{\text{MAP}}(\boldsymbol{\theta}) \mid \boldsymbol{\theta}) \\
&= [\ln \tilde{P}(\xi[m] \mid \boldsymbol{\theta}) - \ln Z(\boldsymbol{\theta})] - [\ln \tilde{P}(\xi^{\text{MAP}}(\boldsymbol{\theta}) \mid \boldsymbol{\theta}) - \ln Z(\boldsymbol{\theta})] \\
&= \ln \tilde{P}(\xi[m] \mid \boldsymbol{\theta}) - \ln \tilde{P}(\xi^{\text{MAP}}(\boldsymbol{\theta}) \mid \boldsymbol{\theta}) \\
&= \sum_i \theta_i [f_i[\xi[m]] - f_i[\xi^{\text{MAP}}(\boldsymbol{\theta})]].
\end{aligned}
$$

- "collapsing" problem

- For conditional density models, we can change the objective to the following, which prevents collapsing

$$\ln P_{\boldsymbol{\theta}}(y[m] \mid x[m]) - \left[\max_{\boldsymbol{y} \neq \boldsymbol{y}[m]} \ln P_{\boldsymbol{\theta}}(y \mid x[m])\right].$$

Find $\gamma, \boldsymbol{\theta}$
that maximize $\gamma$

subject to

$$\ln P_{\boldsymbol{\theta}}(y[m] \mid x[m]) - \ln P_{\boldsymbol{\theta}}(y \mid x[m]) \geq \gamma \qquad \text{for all } m, y \neq y[m]$$

$$\boldsymbol{\theta}^T(f(y[m], x[m]) - f(y, x[m])) \geq \gamma.$$

To prevent margin blowing up we bound \theta

Simple-Max-Margin

Find $\boldsymbol{\theta}$
that minimize $\|\boldsymbol{\theta}\|_2^2$

QP: quad obj+linear constraints

subject to

$$\boldsymbol{\theta}^T(f(y[m], x[m]) - f(y, x[m])) \geq 1 \qquad \text{for all } m, y \neq y[m]$$

31

# Slack variables

- We want to minimize ||w||^2 st

$$\forall_i \forall_{Y_i' \neq Y_i} \log p(Y_i|w, X_i) - \log p(Y_i'|w, X_i) \geq 1.$$

- But we may not be able to achieve this gap, so we introduce slack variables (results in a Hidden Markov Support Vector Machine)

$$\min_{w, \xi} \sum_i \xi_i + \lambda \|w\|_2^2,$$

$$s.t. \quad \forall_i \forall_{Y_i' \neq Y_i} \log p(Y_i|w, X_i) - \log p(Y_i'|w, X_i) \geq 1 - \xi_i, \quad \forall_i \xi_i \geq 0$$

Thanks to Mark Schmidt

# Margin rescaling

- Intuitively if Yi' is similar to Yi, we don't mind if their probabilities are similar, but if they are very different, we want the gap to grow

- This gives max-margin markov network (M3N) aka structural SVM

$$\min_{w,\xi} \sum_i \xi_i + \lambda \|w\|_2^2,$$

$$. \quad \forall_i \forall_{Y_i' \neq Y_i} \log p(Y_i|w, X_i) - \log p(Y_i'|w, X_i) \geq \Delta(Y_i, Y_i') - \xi_i, \quad \forall_i \xi_i \geq 0,$$

Thanks to Mark Schmidt

# Unconstrained form

- We can eliminate the slack vars to get

$$\min_w \sum_i \max_{Y_i' \neq Y_i} (\Delta(Y_i, Y_i') - \log p(Y_i|w, X_i) + \log p(Y_i'|w, X_i))^+ + \lambda ||w||_2^2,$$

- Requires 2$^{nd}$ best decoding. But since $\Delta(Y_i, Y_i) = 0$ we can write

$$\min_w \sum_i \max_{Y_i'} (\Delta(Y_i, Y_i') + \log p(Y_i'|w, X_i)) - \log p(Y_i|w, X_i) + \lambda ||w||_2^2,$$

- This can use generic MAP decoders that just change the local evidence potentials on Y'.

- For associative markov nets, globally optimal.

Thanks to Mark Schmidt

# Cutting plane optimization

- Many possible optimization methods
- Simple approach for QP is cutting planes:
- Maximize quad objective with empty set of constraints – this is an upper bound.
- Add a violated constraint (*)
- Repeat until no violations.
- Thm: only need to add a poly num constraints.
- To find if constraints are violated: define

$$y^{map} = \arg \max_{y \neq y[m]} \tilde{P}(y, x[m]).$$

- If P(y[m],x[m]) < p(ymap,x[m]) +1, add this violation. Else all constraints for m'th case are ok

$$\tilde{P}(y[m], x[m]) > \tilde{P}(y^{map}, x[m]) + 1 \geq \tilde{P}(y, x[m]) + 1,$$