

Max Entropy Priors (17)

Max ignorant priors: prior w/ largest entropy
 $-p^* = \arg\max_{p \in \Delta} H(p) = \arg\max_{p \in \Delta} D(p||U)$ where U is uniform distr.
 $\Rightarrow p^*(y) = \exp(-x^T t(y) - \alpha(x))$ s.t. $E_{p^*}[t(\cdot, y)] = \bar{t}$
Infinite Alphabet case: same as above

Conjugate Priors (18)

Conditionally iid Model: $P_{y_1, \dots, y_N | x} = \prod_{n=1}^N P_{y_n | x}$

Exchangable: (y_1, \dots, y_N) exchangable if \forall perm π ,
 $P_{y_1, \dots, y_N}(y_1, \dots, y_N) \leq P_{y_1, \dots, y_N}(\pi(y_1), \dots, \pi(y_N))$

Theorem: y_1, y_2, \dots is ind. exchangable if $\forall N \exists \bar{x}$,
Cond iid model $P_{y_N | x}$ s.t. $P(y_1, \dots, y_N) = \prod_{n=1}^N P_{y_n | x}$

Conjugate Prior Family: $\mathcal{Q} = \{q(\cdot | \theta) \in \mathcal{P}^X : \theta \in \Theta \subseteq \mathbb{R}^K\}$
s.t. mapping from $\Theta \rightarrow \mathcal{P}^X$ cont. & invertible. \mathcal{Q} is CF Fam if
 $\forall y \in \mathcal{Y}$, we have $p_x(\cdot) \in \mathcal{Q} \Rightarrow p_{x|y}(\cdot | y) \in \mathcal{Q}$.
 \hookrightarrow equivalent to $p_x(\cdot) \in \mathcal{Q} \Rightarrow p_{x|y_1, \dots, y_N}(\cdot | y_1, \dots, y_N) \in \mathcal{Q}$

Natural Conjugate Prior Family: Natural if $\forall y, \exists$ param $\theta(y)$
s.t. $q(\cdot | \theta(y)) \propto p_{y|x}(y | x)$

Theorem 2: For cond. iid models s.t. $\exists \mathbb{R}$ region, $p_{y|x}(\cdot | x)$ cont. func, if conj. prior fam exists, then for every $N \geq 1$ \exists cont. func $t_N(\cdot, \dots, \cdot)$ s.t. $t_N(y_1, \dots, y_N)$ is a sufficient stat wrt inf. abt x , and t_N finite dim ind. of N .

Conjugate Prior of Exp Families: If $p_{y|x} = \exp(\lambda^T t(y) - \alpha(\lambda) + \beta(\lambda))$
then $\mathcal{Q} = \{q(\cdot | t, N) = \exp\{\lambda^T t(x) - N\alpha(\lambda) - \gamma(t, N)\}\}$
is natural conj. prior fam.

Info Geometry of ML & EM (19)

$\hat{x}_{ML}(y) = \arg\max_x p_y(y | x)$
Fact: $\frac{1}{N} \sum_{n=1}^N f(u_n) = E_{\hat{p}_L(y)}[f(u)]$
Fact: $\hat{x}_{ML}(y) = \arg\min_x D(\hat{p}_L(\cdot | y) || p_{y|L}(\cdot | y))$
 \hookrightarrow M-projection (minimizing second term of KL divergence).

Direct EM Algo:

E-step: Compute $U(x; \hat{x}^{(k-1)})$, where
 $U(x; x') = \sum_{z \in \mathcal{Z}} p_{z|x'}(z | y, x') \log p_z(z | x')$
M-step: Compute $\hat{x}^{(k)} = \arg\max_x U(x; \hat{x}^{(k-1)})$.
If $p_{z|x}(z | y, x) = \prod_{n=1}^N p_{z_n | x_n}$, $\hat{x}_n = \arg\max_{x_n} p_{z_n | x_n}$, then E-step simplifies to
 $U(x; x') = \sum_{n=1}^N \sum_{z_n \in \mathcal{Z}_n} p_{z_n | x'}(z_n | y, x_n) \log p_{z_n}(z_n | x_n)$.
Note $p_{z_1 | x}(z_1 | y, x) = \frac{p_z(z_1 | y)}{p_y(y | x)} q_{z_1}(z_1)$, $p_{z_2 | x}(z_2 | y, x) = \sum_{z_1 \in \mathcal{Z}_1} p_z(z_1, z_2 | y, x)$
 \rightarrow want $x = \arg\min_x D(p_{z_1 | x} || p_{z_1 | x'})$.
Let $\hat{p}_z^*(\cdot | x) = \arg\min_{\hat{p}_z \in \hat{\mathcal{P}}^{\mathcal{Z}}_z} D(\hat{p}_z(\cdot | x) || p_{z|L}(\cdot | x))$
 $= \frac{p_z(z | x)}{p_y(y | x)} \hat{p}_z(z | x)$
 \Rightarrow get $\frac{1}{N} U(x; x') = -D(\hat{p}_z^*(\cdot | x) || p_{z|L}(\cdot | x)) - H(\hat{p}_z^*(z | x))$

New EM: E-step: $\hat{p}_z^*(\cdot | x^{(k-1)}) = \arg\min_{\hat{p}_z \in \hat{\mathcal{P}}^{\mathcal{Z}}_z} D(\hat{p}_z(\cdot | x) || p_{z|L}(\cdot | x^{(k-1)}))$
M-step: $\hat{x}^{(k)} = \arg\max_x D(\hat{p}_z^*(\cdot | x^{(k-1)}) || p_{z|L}(\cdot | x))$

Data Processing Inequality (DPI): $g: \mathcal{Z} \rightarrow \mathcal{Z}$, p_{y_1, y_2} induced by p_{z_1, z_2} .
Then $D(p_{y_1} || p_{y_2}) \geq D(p_{y_1|g} || p_{y_2|g})$ w/ equality iff $\frac{p_{z_1}(z_1)}{p_{z_2}(z_2)} = \frac{p_{y_1}(y_1)}{p_{y_2}(y_2)}$
Theorem 1: $p_{y|L}(y | x) = \exp\{\sum_k x_k t_k(y) - \alpha(x) + \beta(x)\}$, then \hat{x}_{ML} satisfies $\frac{1}{N} \sum_{n=1}^N t_k(y_n) = E_{\hat{p}_L(\cdot | \hat{x}_{ML})}(t_k(y)) = E_{\hat{p}_L(\cdot | \hat{x}_{ML})}(t_k(y))$

Typical Seqs, Large Deviations (20)

LLN: If u_1, \dots, u_N iid r.v. w/ mean μ , $E[u_N] < \infty$, then
 $\lim_{N \rightarrow \infty} P\left(\left|\frac{1}{N} \sum_{n=1}^N u_n - \mu\right| > \epsilon\right) = 0$ (conv. in prob.)
Typical Set: Note $\hat{L}_N(y) = \frac{1}{N} \sum_{n=1}^N \log p(y_n) \rightarrow -H(p)$
 \hookrightarrow Defn: seq y is ϵ -typical wrt p if $|\hat{L}_N(y) + H(p)| \leq \epsilon$
 $\mathcal{T}_N^\epsilon(p) = \{y \in \mathcal{Y}^N : |\hat{L}_N(y) + H(p)| \leq \epsilon\}$

Theorem (Asymptotic Equipartition Property): If $\mathcal{T}_N^\epsilon(p)$ is ϵ -typical set, then $\lim_{N \rightarrow \infty} P\{\mathcal{T}_N^\epsilon(p)\} = 1$, $2^{-N(H(p)+\epsilon)} \leq P\{\mathcal{T}_N^\epsilon(p)\} \leq 2^{-N(H(p)-\epsilon)}$
and $(1-\epsilon)2^{-N(H(p)+\epsilon)} \leq |\mathcal{T}_N^\epsilon(p)| \leq 2^{N(H(p)+\epsilon)}$ for large enuf N .
Divergence ϵ -typical wrt p relative to q : $q: |\hat{L}_N(y) - D(p||q)| \leq \epsilon$
 $\hookrightarrow \mathcal{T}_N^\epsilon(p||q)$
 $\lim_{N \rightarrow \infty} P\{\mathcal{T}_N^\epsilon(p||q)\} = 1$
Theorem: $(1-\epsilon)2^{-N(D(p||q)+\epsilon)} \leq Q\{\mathcal{T}_N^\epsilon(p||q)\} \leq 2^{-N(D(p||q)-\epsilon)}$

Cramer's Theorem: If $y = (y_1, \dots, y_N)$ iid from q , $t: \mathbb{R}^d \rightarrow \mathbb{R}$ s.t. $\mu = E_q[t(y)] < \infty$, then for any γ, ρ , have:
 $\lim_{N \rightarrow \infty} -\frac{1}{N} \log P\left(\frac{1}{N} \sum_{n=1}^N t(y_n) \geq \gamma\right) = E_q[t]$
where $E_q[t] \triangleq D(p_L(\cdot | x) || q)$ w/ $p_L(x) = q(y) e^{x^T t(y) - \alpha(x)}$,
and x s.t. $E_{p_L}[\cdot](t(y)) = \gamma$.
Furthermore, we can express $E_q[t] = D(q || p_L)$ w/ $p_L(x) = q(y) e^{x^T t(y) - \alpha(x)}$
Can calculate: $\frac{1}{\log e} \cdot \frac{1}{\partial_x} (D(p_L(\cdot | x) || q)) = x^T t(y) > 0 \quad \forall x > 0$

Method of Types + Sanov (22)

Type: prob distr. $\hat{p}(y) = \frac{1}{N} \sum_{n=1}^N 1_{y=y_n} = \frac{N_y(y)}{N}$
Set of Types: $\mathcal{P}_N^X =$ set of all possible types
Type class: $J_N^x(p) = \{y \in \mathcal{Y}^N : \hat{p}(y) = p\}$
Exponential Rate Notation: $f(N) = 2^{N\alpha}$ means $\lim_{N \rightarrow \infty} \frac{\log f(N)}{N} = \alpha$.
Properties: $|J_N^x(p)| \leq (N+1)^{|\mathcal{Y}|}$
Lemma: If $z \in \mathcal{P}^X$ defined over \mathcal{Y} , then $q^z(y) = 2^{-N D(z || q)}$ (prob and self info approx)
Lemma: $|J_N^x(p)| \leq 2^{N H(p)}$, $c N^{-|\mathcal{Y}|} 2^{N H(p)} \leq |J_N^x(p)| \leq 2^{N H(p)}$
Theorem 1: $c N^{-|\mathcal{Y}|} 2^{-N D(p||q)} \leq Q\{\mathcal{T}_N^x(p)\} \leq 2^{-N D(p||q)}$
Sanov's Theorem: $S \subset \mathcal{P}^X$ arbitrary, $q \in \mathcal{P}^X$ arbitrary,
then $Q\{S \cap \mathcal{P}_N^x\} \leq (N+1)^{|\mathcal{Y}|} 2^{-N D(p||q)} = 2^{-N D(p||q)}$
where $p^* = \arg\min_{p \in S} D(p||q) = \mathcal{I}\text{-proj of } q \text{ onto } S$.
If $d(S) = d(\text{int}(S))$, then $Q\{S \cap \mathcal{P}_N^x\} = 2^{-N D(p^*||q)}$
Pf sketch: bound size $|S|$ w/ $(N+1)^{|\mathcal{Y}|}$ and $Q\{S \cap \mathcal{P}_N^x\} \leq c N^{-|\mathcal{Y}|} 2^{-N D(p^*||q)}$

Conditional Limit Thm: If $S \subset \mathcal{P}^X$ closed & convex, $y \in \mathcal{Y}^N$ iid wrt $q \in \mathcal{P}^X$,
then for any $\epsilon > 0$ we have $\lim_{N \rightarrow \infty} P\left(\left|\frac{1}{N} \sum_{n=1}^N t(y_n) - p_0\right| > \epsilon \mid \frac{1}{N} \sum_{n=1}^N t(y_n) \in S\right) = 0$
where $p_0 = \mathcal{I}\text{-proj from } q \text{ to } S$ (the type $\rightarrow p_0$)
Pinsker ineq: For any $P, Q \in \mathcal{P}^X$, we have $\|P - Q\|_1 = \sum_{y \in \mathcal{Y}} |p(y) - q(y)| \leq \sqrt{2 \log 2 D(P||Q)}$

Asymptotics of Hypo Testing (23)

Setup: Binary hypo test $H_0: p_1(y) = \prod_{n=1}^N p_{1n}(y_n)$, $H_1: p_2(y) = \prod_{n=1}^N p_{2n}(y_n)$
 \rightarrow log-likelihood ratio: $\hat{L}(y) = \frac{1}{N} \sum_{n=1}^N \log \frac{p_{1n}(y_n)}{p_{2n}(y_n)}$
Typicality based rule:
 $\hat{L}(y) \geq \tau \Rightarrow p_1$
 $\hat{L}(y) < \tau \Rightarrow p_2$
where $\tau = \frac{1}{2} \log \frac{p_1(p_1)}{p_2(p_2)}$
 \hookrightarrow get $P_F \leq \epsilon + 2^{-N D(p_1||p_2)}$, $P_M \leq \epsilon + 2^{-N D(p_2||p_1)}$

Theorem: For log-likelihood ratio test w/ threshold τ ($\hat{L}(y) \geq \tau \Rightarrow p_1$),
with $\hat{L}(y) = \frac{1}{N} \sum_{n=1}^N \log \frac{p_{1n}(y_n)}{p_{2n}(y_n)}$, we have:
 $P_F = 2^{-N D(p_1||p_2)}$, $P_M = 2^{-N D(p_2||p_1)}$ where $\hat{p}_1 = H_1$
 $P_F \approx \frac{1}{2} p_1^{1-x} p_2^x$ w/ x s.t. $D(p_x || p_2) = D(p_1 || p_2) = \tau$.

NP-hypo testing: Constrain P_F , minimize P_M
 \rightarrow If P_F doesn't decay exponentially, let $\tau = -D(p_1 || p_2)$, $P_F = p_1 - P_M = 2^{-N D(p_1 || p_2)}$

Bayesian Hypo Testing: $E[C] = C_0 p_0 P_F + C_1 p_1 P_M = 2^{-N \min(D(p_1 || p_2), D(p_2 || p_1))}$
 \rightarrow satisfied at $\delta = 0$, $E_C = D(p_1 || p_2) \approx D(p_2 || p_1)$

Asymptotics of Param Est. (25)

Identifiability: $x \neq x'$ s.t. $D(p_{y|L}(\cdot | x) || p_{y|L}(\cdot | x')) > 0 \Leftrightarrow x \neq x'$
Consistency: $\hat{x}_N(y^N) \xrightarrow{P} x$ as $N \rightarrow \infty$ (weak) / $\xrightarrow{a.s.}$ (strong)
Asymptotic Normality: $\sqrt{N}(\hat{x}_N - x) \xrightarrow{d} N(0, \sigma^2(x))$ for $\sigma^2(x) > 0$
Asymptotic Efficiency: Asymptotic Normal and $\sigma^2(x) = 1/J_2(x)$ (Fisher info)
Theorem 1: If $|\hat{L}_1|, |\hat{L}_2| < \infty$, y_1, y_2 iid $\sim p_L(\cdot | x)$ for identifiable x ,
and $p_{y_1|y_2} > 0 \quad \forall y_1, y_2$, then ML estimate $\hat{x}_N(y^N)$ weakly consistent.
Theorem 2: If $p_{y_1|y_2} = \exp(-\alpha(y_1) - \alpha(y_2) + \beta(y_1, y_2))$, y^N iid $\sim p_L(\cdot | x)$ s.t. $\hat{\alpha}(x) = \bar{\alpha}$, $\hat{\beta}(x) = \bar{\beta}$,
then $\hat{x}_N(y^N)$ is strongly consistent.
Theorem 3: same as above, now $\hat{x}_N(y^N)$ is asymptotically efficient.
Folk thm: $\hat{x}_N(y^N) \sim N(x, \frac{1}{N J_2(x)})$

Theorem 6: If $P = \{p_{y|L}(\cdot | x) : x \in \mathcal{X}, |\mathcal{X}|, |\mathcal{Y}| < \infty\}$,
 y_1, y_2 iid according to $p_{y|L} \notin P$, $p_{y_1|y_2} > 0$, $p_{y|L}(y | x) > 0 \quad \forall y, x$,
and $x_0 = \arg\min_{x \in \mathcal{X}} D(p_{y_1|y_2} || p_{y|L}(\cdot | x))$ unique. Then $\hat{x}_N(y_1, y_2) \xrightarrow{P} x_0$.

Laplace Method: Approximate $\int_{\mathcal{X}} e^{N g(x)} dx$ if g twice diff, unique max at $x_0 \in \mathcal{X}$, s.t. $\hat{g}(x_0) < 0$, then
 $\frac{\int_{\mathcal{X}} e^{N g(x)} dx}{e^{N g(x_0)} \sqrt{\frac{2\pi}{N |\hat{g}''(x_0)|}}} \rightarrow 1$ as $N \rightarrow \infty$
We want to approximate posterior $p_{x|y_N}(x | y^N) = \frac{p_x(x) \prod_{n=1}^N p_{y_n|x}(y_n | x)}{\int p_x(x) \prod_{n=1}^N p_{y_n|x}(y_n | x) dx}$
If \mathcal{X} cont. alphabet, then $p_{x|y_N}(x | y^N) \sim N(\hat{x}_N(y^N), \frac{1}{N J_2(x)})$.

Theorem 8: If $\tilde{x} = \sqrt{N}(x - \hat{x}_N(y^N))$, y^N iid from $p_{y_1|y_2}(y_1, y_2)$, $x \in \mathcal{X}$,
 $\hat{\alpha}(x) = \bar{\alpha}$, and $p_{y_1|y_2} = \exp(-\alpha(y_1) - \alpha(y_2) + \beta(y_1, y_2))$,
then $\ln(p_{x|y_N}(\tilde{x} | y^N) / p_{x|y_N}(0 | y^N)) \xrightarrow{P} -\frac{1}{2} \tilde{x}^T J_2(x) \tilde{x}$ as $N \rightarrow \infty$.

Asymptotics of Inf. & Universality (26)

Universal inference: After enough data, can make predictions abt future data.
Defn Universal Predictor: $S_N = \{p_{y|L}(\cdot | x) \in \mathcal{P}^X, x \in \mathcal{X}\}$ family of models
 $q_1, q_2, \dots, q_N \in \mathcal{P}^X$ seq. of predictors, $q_N(\cdot | y^{N-1})$ est. of dist given y^{N-1} ,
then seq. is universal if $\lim_{N \rightarrow \infty} \frac{1}{N} \log \frac{p_{y|L}(y^N)}{q_N(y^N)} = 0$ where
 $\bar{p}_N(x) = \frac{1}{N} \sum_{n=1}^N p_n(x)$ w/ $p_n(x) = E_{p_{y|L}(\cdot | x)} \left[\log \frac{p_{y_n|y^{n-1}}(y_n | y^{n-1})}{q_n(y_n | y^{n-1})} \right]$
and $p_{y_n|y^{n-1}}(y_n | y^{n-1}) = \frac{p_{y_n}(y_n)}{p_{y^{n-1}}(y^{n-1})}$.

\hookrightarrow A cor: $p_n(x) = E_{p_{y|L}(\cdot | x)} [D(p_{y_n|y^{n-1}}(\cdot | y^{n-1}) || q_n(\cdot | y^{n-1}))]$
 $\hookrightarrow \bar{p}_N(x) = \frac{1}{N} D(p_{y|L}(\cdot | x) || q_{y^N}(\cdot))$, where $q_{y^N}(\cdot) = \prod_{n=1}^N q_n(y_n | y^{n-1})$
Theorem 2: If model capacity of model fam S_N is C_N ,
 $y^N = (y_1, \dots, y_N)$ from $p \in S_N$, then universal if possible iff $\lim_{N \rightarrow \infty} \frac{C_N}{N} = 0$.

Asymptotic Least informative prior: $p_x^* \in \mathcal{P}^X$ for fam w/
capacity C_N if $I(x; y^N)$ associated w/ $p_{y^N}(y^N | x) p_x^*(x)$
satisfies $(n - I(x; y^N)) = o(1)$ as $N \rightarrow \infty$.
 \hookrightarrow Construct prediction $q_{y_1|y^{n-1}}(y_n | y^{n-1}) = \frac{q(y^N)}{q(y^{n-1})}$ w/ p^* prior.
Theorem 2: Model fam \mathcal{S} , $|\mathcal{X}|, |\mathcal{Y}| < \infty$, $y^N \sim p_{y|L}(\cdot | x)$
where $x \in \mathcal{X}$, $p_L(\cdot | x) > 0 \quad \forall y$, then:
 $D(p_{y^N|L}(\cdot | x) || p_{y^N}) = -\log p_L(x) + o(1)$ and $I(x; y^N) = H(p_L) + o(1)$
and $C_N = \log |\mathcal{X}| + o(1)$ as $N \rightarrow \infty$, achieved by uniform prior $p_x^* = \frac{1}{|\mathcal{X}|}$.

Theorem 4: $\mathcal{S} = \{ \exp(-\alpha(y) - \alpha(y_1) + \beta(y, y_1)) \}$, $y^N \sim p_{y|L}(\cdot | x)$, then
 $D(p_{y^N|L}(\cdot | x) || p_{y^N}) = \frac{1}{2} \log \frac{N J_2(x)}{2\pi e} - \log p_L(x) + o(N)$, $N \rightarrow \infty$.
Cor: If $\int_{\mathcal{X}} \sqrt{J_2(x)} dx < \infty$ then $I(x; y^N) = \frac{1}{2} \log \frac{N}{2\pi e} - D(p_L || p_x^*) + \log \left(\int_{\mathcal{X}} \sqrt{J_2(x)} dx \right) + o(1)$
where $p_x^* = \sqrt{J_2(x)} / \int_{\mathcal{X}} \sqrt{J_2(x)} dx \leftarrow$ Jeffreys' prior

Theorem 3: For $\mathcal{S} = \{ \exp(-\alpha(y) - \alpha(y_1) + \beta(y, y_1)) \}$, $\tilde{x} = \sqrt{N}(x - \hat{x}_N(y^N))$,
we have $D(p_{x|y^N}(\tilde{x} | y^N) || N(0, 1/J_2(x))) \xrightarrow{P} 0$ achieved model capacity

Theorem 7: If $y^N \in \mathcal{Y}^N$ generated iid wrt $p_L \notin P$, then no predictor is universal. Asymptotic loss $\bar{p}_N(y^N) \rightarrow \frac{1}{N} D(p_{y^N|L} || p_{y^N}) \geq \rho_{\min} > 0$
where $q_{y^N}(y^N) = \prod_{n=1}^N q_n(y_n)$, $\rho_{\min} = \min_{x \in \mathcal{X}} D(p_{y|L} || p_{y|L}(\cdot | x))$.

Tradeoff between class size - large $\Rightarrow C_N$ large, small \Rightarrow not in class