

## Binary hypothesis testing (2)

Find  $H \in \{H_0, H_1\}$  prior  $P_0 = \mathbb{P}[H = H_0]$  and  $P_1 = 1 - P_0$ . Model consists of  $p_H, p_{Y|H} \rightarrow \hat{H}(y)$ . With cost function  $C_{ij}$  (guess  $i$ , truth  $j$ ), we seek  $\hat{H} = \operatorname{argmin}_{f(\cdot)} \varphi(f)$ , for  $\varphi(f) = \mathbb{E}[C(\hat{H}, f(y))]$ . The **likelihood ratio test** minimizes **Bayes' risk**  $\varphi$ :

$$L(y) := \frac{p_{Y|H}(y|H_1)}{p_{Y|H}(y|H_0)} \underset{H_0}{\overset{H_1}{\geq}} \frac{P_0(C_{10} - C_{00})}{P_1(C_{01} - C_{11})} := \eta$$

Proof sketch:  $f(y)$  can be computed pointwise on  $y$ , so compare the expected cost for  $\hat{H}(y) = H_0, H_1$ .

- **maximum a-posteriori rule**:  $C$  is 0-1 loss, rule compare  $p_{H|Y}(H_i|y)$ .
- **maximum likelihood rule**: if also  $P_1 = P_2 = \frac{1}{2}$ , rule compare  $p_{Y|H}(y|H_i)$ .

## Neyman-Pearson/OCs (3,4)

$P_F = \mathbb{P}[\hat{H} = H_1 | H = H_0]$ , and  $P_D = \mathbb{P}[\hat{H} = H_1 | H = H_1]$ . For LRTs  $L(y) \underset{\eta}{\geq} \eta$ , the **operating characteristic** is set  $\{(P_F, P_D) \mid \eta\}$ . The OC is  $(1, 1) \rightarrow (0, 0)$  as  $\eta$  decreases, and is convex. Bayesian framework selects  $\eta : (P_F, P_D)$  on LRT minimizing  $\alpha P_F - \beta P_D + \gamma$  (determined by costs).

$$\zeta_{NP}(\alpha) = \operatorname{argmax}_{\hat{H}(\cdot)} P_D \text{ s.t. } P_F \leq \alpha.$$

$\zeta_{NP}(\alpha)$  is concave nondecreasing from  $(0, 0)$  to  $(1, 1)$ .

*Continuous case*: if  $L(y)$  continuous then the  $\operatorname{argmax} \hat{H}(\cdot)$  is of the form  $\mathbf{1}[L(y) \geq \eta]$ .

*Discrete case*: equality in the LRT is possible; randomization  $p_{\hat{H}|Y}(\cdot|y)$  can improve over deterministic decision rules. In particular, for rule  $\hat{H}$  that picks  $\hat{H}'$  with probability  $p$  and  $\hat{H}''$  with probability  $1 - p$ , we have  $P_D(\hat{H}) = pP_D(\hat{H}') + (1 - p)P_D(\hat{H}'')$ , and same for  $P_F$ .

*Neyman-Pearson Lemma*: There exists an optimum Neyman-Pearson rule of the form  $q_*(y) = 0$  if  $L(y) < \eta$ ,  $p$  if  $L(y) = \eta$ , and 1 if  $L(y) > \eta$ .

## Minimax hypothesis testing (5)

Costs  $C_{ij}$  known, but nature picks worst prior. Let  $\varphi(p, r)$  be Bayes risk of  $r(\cdot) = p_{\hat{H}|Y}(\cdot|y)$  under prior  $p$ . Seek  $r_M(\cdot) = \operatorname{argmin}_{r(\cdot)} \max_p \varphi(p, r)$ .

The **mismatched Bayes risk**  $\varphi_B(p, q, \lambda)$  is risk  $\varphi(p, r_B(\cdot; q, \lambda))$ , where  $r_B(\cdot; q, \lambda)$  is the Bayes rule for prior  $q$ . Hence **Bayes risk** is  $\varphi_B^*(p) = \varphi_B(p, p, \lambda)$ , independent of  $\lambda$ . Note  $\varphi_B$  is linear in  $p$ , and minimized over  $q$  at  $q = p$ . Also  $\varphi_B^*(p)$  is concave, and  $\varphi_B^*(0) = C_{00}, \varphi_B^*(1) = C_{11}$ .

The **minimax decision rule** is  $r_B(\cdot; p_*, \lambda_*)$ , where  $p_*, \lambda_*$  correspond to the point  $(P_F^*, P_D^*)$  at the intersection of  $\zeta_{NP}(P_F^*)$  and

$$g_M(P_F^*) := \frac{C_{01} - C_{00}}{C_{01} - C_{11}} - \frac{C_{10} - C_{00}}{C_{01} - C_{11}} P_F^*.$$

If there is no intersection, set  $p_* = 0$  if  $\zeta_{NP} > g_M(P_F)$  and  $p_* = 1$  otherwise, and  $\lambda$  can be set arbitrarily

**Minimax Inequality**: in general,

$$\min_a \max_b g(a, b) \geq \max_b \min_a g(a, b).$$

In this case,

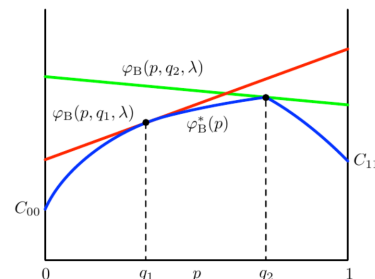
$$\min_{r(\cdot)} \max_p \varphi(p, r) = \max_p \min_{r(\cdot)} \varphi(p, r)$$

by Von Neumann's Theorem; one can show  $\varphi(p, r)$  has a saddle point.

A **equalizer rule** is one such that

$$\mathbb{E}[C(H, \hat{H}) | H_0] = \mathbb{E}[C(H, \hat{H}) | H_1],$$

i.e. the Bayes risk of the rule is the same regardless of the prior  $p$ . The minimax rule is equalizer rule if  $p_* \in \{0, 1\}$ .



## Bayesian parameter estimation (6)

With the posterior  $p_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) = \frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})}{\int p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})p_{\mathbf{x}}(\mathbf{x})d\mathbf{x}}$ , choose  $\hat{\mathbf{x}} = \operatorname{argmin}_{f(\cdot)}[C(\mathbf{x}, f(\mathbf{y}))]$  for some cost function  $C$ .  $\hat{\mathbf{x}}(\mathbf{y})$  should be minimized pointwise.

Let error  $e(\mathbf{x}, \mathbf{y}) = \hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}$ , so bias  $b = \mathbb{E}[e(\mathbf{x}, \mathbf{y})]$ . **Error covariance matrix**  $\Lambda_e = \mathbb{E}[(e-b)(e-b)^\top]$ , **error correlation matrix**  $\mathbb{E}[ee^\top] = \Lambda_e + bb^\top$ . MSE is  $\operatorname{tr}(\mathbb{E}[ee^\top])$ .

- **min absolute error**:  $C(x, \hat{x}) = |x - \hat{x}|$  gives the median of the posterior.
- **min uniform cost**  $C(x, \hat{x}) = \mathbf{1}[|x - \hat{x}| > \varepsilon]$  gives the mode as  $\varepsilon \rightarrow 0$  (MAP estimator).

**Bayesian Least Squares**  $C(a, \hat{a}) = \|a - \hat{a}\|^2$  gives mean of posterior,  $\hat{\mathbf{x}}(\mathbf{y}) = \mathbb{E}[\mathbf{x}|\mathbf{y} = \mathbf{y}]$ . The BLS estimator is unbiased, with error covariance & correlation matrix  $\mathbb{E}[\Lambda_{\mathbf{x}|\mathbf{y}}(\mathbf{y})]$ .

*Orthogonality*:  $\hat{\mathbf{x}}$  is  $\hat{\mathbf{x}}_{\text{BLS}}$  iff unbiased and

$$\mathbb{E}[(\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x})g(\mathbf{y})^\top] = 0 \text{ for all } g(\cdot).$$

## Linear least-squares\* (7)

**LLS estimator**  $\hat{\mathbf{x}}_{\text{LLS}}(\mathbf{y}) = \operatorname{argmin}_{f(\cdot) \in \mathcal{B}} \mathbb{E}[\|\mathbf{x} - f(\mathbf{y})\|^2]$ , where  $\mathcal{B} = \{f(\mathbf{y}) = \mathbf{A}\mathbf{y} + \mathbf{d}\}$ .

*Orthogonality*:  $\hat{\mathbf{x}}$  is the LLS estimator iff unbiased and  $\mathbb{E}[(\hat{\mathbf{x}} - \mathbf{x})\mathbf{y}^\top] = 0$ . This is equivalent to  $\mathbb{E}[(\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x})(\mathbf{F}\mathbf{y} + \mathbf{g})^\top] = 0 \quad \forall \mathbf{F}, \mathbf{g}$ . Closed-form:

$$\hat{\mathbf{x}}_{\text{LLS}}(\mathbf{y}) = \mu_{\mathbf{x}} + \Lambda_{\mathbf{xy}}\Lambda_{\mathbf{y}}^{-1}(\mathbf{y} - \mu_{\mathbf{y}}).$$

**Gaussian** random variables  $x = \frac{1}{\sqrt{2\pi\sigma^2}}e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ , vectors  $p_{\mathbf{x}}(\mathbf{x}) = \frac{1}{\sqrt{|2\pi\Lambda|}}e^{-\frac{1}{2}(\mathbf{x}-\mu)^\top\Lambda^{-1}(\mathbf{x}-\mu)}$ . Then  $a^\top\mathbf{x}$  is Gaussian. If  $\mathbf{x}, \mathbf{y}$  are jointly Gaussian,  $\hat{\mathbf{x}}_{\text{BLS}} = \hat{\mathbf{x}}_{\text{LLS}}$ .

## Nonbayesian param. estimation (8)

Only consider *valid*  $f(\cdot)$ , i.e. independent of  $x$ . Let  $e(\mathbf{y}) = \hat{\mathbf{x}} - \mathbf{x}$ ,  $b_{\hat{\mathbf{x}}}(x) = \mathbb{E}[e(\mathbf{y})]$ , and  $\Lambda_e(x) = \mathbb{E}[(e-b)(e-b)^\top]$ . MSE is trace of

$$\mathbb{E}[e(\mathbf{y})e(\mathbf{y})^\top] = \Lambda_e(x) + b_{\hat{\mathbf{x}}}(x)b_{\hat{\mathbf{x}}}(x)^\top.$$

A **minimum variance unbiased** estimator satisfies  $\lambda_{\hat{\mathbf{x}}}^*(x) \leq \lambda_{\hat{\mathbf{x}}}(x)$  for every  $x$ . This does not have to exist.

**Cramér-Rao**: if  $\mathbb{E}[\frac{\partial}{\partial x} \ln p_{\mathbf{y}}(\mathbf{y}; x)] = 0$  for all  $x$ , then variance of unbiased estimator  $\hat{x}$  is at least

$$\lambda_{\hat{x}}(x) \geq 1/J_{\mathbf{y}}(x) \quad J_{\mathbf{y}}(x) = \mathbb{E}[S(\mathbf{y}; x)^2]$$

$$S(\mathbf{y}; x) = \frac{\partial}{\partial x} \ln p_{\mathbf{y}}(\mathbf{y}; x).$$

**Vector Cramér-Rao**: if  $\mathbb{E}[\frac{\partial}{\partial \mathbf{x}} \ln p_{\mathbf{y}}(\mathbf{y}; \mathbf{x})] = 0$  then for any unbiased  $\hat{\mathbf{x}}(\cdot)$ ,  $\Lambda_{\hat{\mathbf{x}}} = \operatorname{Cov}[\hat{\mathbf{x}}|\mathbf{x}]$  satisfies

$$\Lambda_{\hat{\mathbf{x}}} - J_{\mathbf{y}}^{-1}(\mathbf{x}) \succeq 0, \quad J_{\mathbf{y}}(\mathbf{x}) = \mathbb{E}[S(\mathbf{y}; \mathbf{x})S(\mathbf{y}; \mathbf{x})^\top]$$

$$S(\mathbf{y}; \mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} \ln p_{\mathbf{y}}(\mathbf{y}; \mathbf{x}).$$

Also,  $\mathbb{E}\left[\left(\frac{\partial}{\partial x} \ln p_{\mathbf{y}}(\mathbf{y}; x)\right)^2 + \frac{\partial^2}{\partial x^2} \ln p_{\mathbf{y}}(\mathbf{y}; x)\right] = 0$ , implies  $J_{\mathbf{y}}(x) = -\mathbb{E}\left[\frac{\partial^2}{\partial x^2} \ln p_{\mathbf{y}}(\mathbf{y}; x)\right]$ .

**Efficient estimators** match equality of Cramér-Rao bound for all  $x$ , i.e. is of the form

$$\hat{x}(y) = x + \frac{1}{J_{\mathbf{y}}(x)}S(\mathbf{y}; x),$$

(where the RHS must be independent of  $x$ ). Efficient estimators are unbiased & unique & MVU if existent.

The **Maximum Likelihood estimator** is

$$\hat{x}_{\text{ML}}(y) = \operatorname{argmax}_x p_{\mathbf{y}}(y; x).$$

$\hat{x}_{\text{ML}}$  commutes with invertible transformations: if  $\theta = g(x)$ ,  $\theta_{\text{ML}}(\cdot) = g(\hat{x}_{\text{ML}}(\cdot))$ . Also, if both exist,  $\hat{x}_{\text{ML}}(\cdot) = \hat{x}_{\text{eff}}(\cdot)$ .

**Gauss-Markov Theorem**: If  $\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{w}$  for  $\mathbf{w} \sim N(0, \Lambda)$ , then MVU estimator  $\hat{\mathbf{x}}_{\text{ML}}(\mathbf{y})$  is

$$(\mathbf{H}^\top\Lambda^{-1}\mathbf{H})^{-1}\mathbf{H}^\top\Lambda^{-1}\mathbf{y} = \operatorname{argmin}_{\mathbf{x}} \|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2.$$

## Exponential Families (9)

An **single-param. exponential family** is of form

$$p_{\mathbf{y}}(\mathbf{y}; x) = \exp(\lambda(x)t(\mathbf{y}) - \alpha(x) + \beta(\mathbf{y})),$$

“Natural parameter”  $\lambda(\cdot)$ , “natural statistic”  $t(\cdot)$ , “log base function”  $\beta(\cdot)$ . Note  $\alpha(\cdot)$  normalizes, and  $t(\cdot)$ ,  $\beta(\cdot)$  can be shifted by constants. Call  $q(y) \propto \exp(\beta(y))$  the **base distribution**.

$$\dot{\alpha}(x) = \dot{\lambda}(x) \mathbb{E}[t(y)]$$

$$\ddot{\alpha}(x) = \ddot{\lambda}(x) \mathbb{E}[t(y)] + \dot{\lambda}(x)^2 \mathbb{V}[t(y)]$$

$$J_y(x) = \dot{\lambda}(x) \frac{d}{dx} \mathbb{E}[t(y)]$$

A **canonical exponential family** has  $\lambda(x) = x$ .  $e^{\alpha(x)} = \mathbb{E}_q[e^{xy}]$ ,  $\dot{\alpha}(x) = \mathbb{E}[t(y)]$ ,  $\ddot{\alpha}(x) = \mathbb{V}[t(y)]$ . Examples include the *geometric mean*  $p(\cdot; x) \propto q_1(\cdot)^x q_2(\cdot)^y$  and *tilted*  $p(\cdot; x) \propto q(y) e^{xt(y)}$ .

**Multi-parameter exponential family:**

$$p_y(\mathbf{y}; \mathbf{x}) = \exp(\lambda(x)^\top \mathbf{t}(\mathbf{y}) - \alpha(\mathbf{x}) + \beta(\mathbf{y})).$$

Over finite alphabet  $\mathcal{Y}$ , any distribution  $q$  can be generated by one exponential family  $p_y(y; \mathbf{x}) = \exp(\sum_{i \in \mathcal{Y}} x_i t_i(y) - \alpha(\mathbf{x}))$  via  $x_i = \ln q(i)$ .

## Sufficient Statistics (10)

A *sufficient* statistic  $t(\cdot)$  w.r.t. family  $\{p_y(\cdot; x)\}$  is such that  $p_{y|t}(\cdot; x)$  is independent of  $x$ . Equivalently  $\frac{L_y(x)}{L_t(x)} := \frac{p_y(y; x)}{p_t(t(y); x)}$  is not a function of  $x$ .

**Neyman Factorization:**  $t(\cdot)$  is sufficient iff exist  $a(\cdot, \cdot)$  and  $b(\cdot)$  such that

$$p_y(y; x) = a(t(y), x) b(y).$$

Sufficient statistic  $s(\cdot)$  is **minimal** if for any sufficient  $t$ , exists  $g$  such that  $s = g \circ t$ .

For example, for finite  $\mathcal{X}$ ,  $t(y) = \langle p_y(y; x_i) \rangle$  is sufficient, and  $t(y) = \frac{p_y(y; x_i)}{p_y(y; x_0)}$  is minimal.

A sufficient statistic  $t$  is **complete** if the only  $\varphi$  such that  $\mathbb{E}[\varphi(t(y))] = 0$  for all  $x \in \mathcal{X}$  is  $\varphi(\cdot) \equiv 0$ . Completeness implies minimality: if  $t$  is complete and  $s$  is minimal, let  $s = g(t)$ ; then  $\mathbb{E}[t|s = s]$  is a function of  $s$  and hence  $t$ , so  $t - \mathbb{E}[t|s = s]$  has mean zero, hence is zero, so  $t$  is a function of  $s$ .

**Bayesian formulation:**  $t(\cdot)$  is sufficient w.r.t.  $p_{x,y}$  if  $p_{x|y}(\cdot|y) = p_{x|t}(\cdot|t(y))$ . Neyman factorization looks like  $p_{y|x}(y|x) = p_{t|x}(t(y)|x) p_{y|t}(y|t(y))$ .

*Statistics are partitions:* If  $L_{y_1}(x) \propto L_{y_2}(x)$  then  $y_1, y_2$  provide “the same information” about  $x$ . Sufficient statistics group all  $y_i$  with proportional  $L_{y_i}(\cdot)$  together; minimum sufficient statistics in addition group all  $y_i$  with non-proportional  $L_{y_i}(\cdot)$  apart.

## EM algorithm (12)

Observed data  $y$  from r.v.  $y \sim p_y(\cdot; x)$ . We wish to compute  $\hat{x}_{ML}$ . Find some **complete data**  $\mathbf{z}$  such that  $y = g(\mathbf{z})$ , such that ML estimator of  $\mathbf{z}$  is easy.

1. Set  $\ell = 0$  and initialize  $\hat{x}$  arbitrarily.
2. Find  $U(x; \hat{x}) = \mathbb{E}_{p_{z|y}(\cdot|y; \hat{x})}[\log p_z(\mathbf{z}|x)]$ .
3. Set new  $\hat{x}$  to  $\arg\max_x U(x; \hat{x})$ .

It is possible to show

$$\ell_y(x; y) = U(x, x') - U(x', x') + \ell_y(x'; y),$$

So the new  $\hat{x}$  increases log-likelihood  $\ell_y(\hat{x}; y)$  at each step. This converges to a stationary point.

For *Gaussian mixture models*, i.e. of the form  $p_y(\mathbf{y}; \theta) = \sum_k \pi_k N(\mathbf{y}; \mu_k, \Lambda_k)$ , and complete data is  $\mathbf{y}$  along with the weights  $\pi_i$ . More generally, for *exponential families*, i.e. complete data  $\mathbf{z}$  is in

$$p_z(\mathbf{z}; \mathbf{x}) = \exp(\mathbf{x}^\top \mathbf{t}(\mathbf{z}) - \alpha(\mathbf{x}) + \beta(\mathbf{z})),$$

we have  $U(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbb{E}_{p_{z|y}(\cdot|y; \mathbf{x}')}[\mathbf{t}(\mathbf{z})] - \alpha(\mathbf{x}) + \mathbb{E}_{p_{z|y}(\cdot|y; \mathbf{x}')}[\beta(\mathbf{z})]$ ; solve partials and substitute  $\frac{\partial}{\partial x_i} \alpha(\mathbf{x}) = \mathbb{E}_{p_z(\cdot; \mathbf{x})}[t_i(\mathbf{z})]$  to get new  $\hat{\mathbf{x}}'$  from  $\hat{\mathbf{x}}$ :

$$\mathbb{E}_{p_z(\cdot; \hat{\mathbf{x}}')}[\mathbf{t}_k(\mathbf{z})] = \mathbb{E}_{p_{z|y}(\cdot|y; \hat{\mathbf{x}})}[\mathbf{t}_k(\mathbf{z})].$$

## Decision theory (13)

Inference predicts a probability distribution  $q(\cdot)$  for  $x$ . The only smooth cost function  $C(x, q)$  s.t. **proper**:  $p_{x|y}(\cdot|y) = \arg\min_q \mathbb{E}[C(x, q)|y = y]$  and **local**:  $C(x, q)$  is a function of  $x$  and  $q(x)$ , is **log-loss**:  $-A \log q(x) + B(x)$ , as long as  $|\mathcal{X}| \geq 3$ . **Entropy**:  $H = -\sum_a p(a) \log p(a) = \mathbb{E}[C(x, p_x)]$ . This is concave in  $p$ . **Conditional entropy**:

$H(x|y) = \mathbb{E}[C(x, p_{x|y})]$ ; expands

$$H(x|y) = \sum_y p_y(y) H(x|y=y).$$

Note  $H(x|y) \leq H(x) \leq \log |\mathcal{X}|$ , and  $H(x, y) = H(x|y) + H(y)$ .

**Mutual information**  $I(x; y) = H(x) - H(x|y)$ . This is *symmetric*, since it can be computed as  $\sum_{x,y} p_{x,y}(x, y) \log \frac{p_{x,y}(x, y)}{p_x(x)p_y(y)}$ . Some identities:

- $I(x; y) = 0$  iff  $x, y$  independent
- $I(x; y, z) = I(x; z) + I(x; y|z)$  where  $I(x; y|z) = H(x|z) - H(x|y, z)$
- $H(x, y) = H(x) + H(y) - I(x; y)$ .

**Data Processing Inequality:**  $I(x; y) \geq I(x; t)$ , with equality if and only if  $t(y)$  is sufficient.

**Information (KL) Divergence:** given  $x \sim p(\cdot)$ ,  $D(p||q) = \mathbb{E}[C(x, q) - C(x, p)] = \sum_a p(a) \log \frac{p(a)}{q(a)}$ .  $D(p||U) = \log |\mathcal{X}| - H(p)$ ,  $I(x; y) = D(p_{x,y}||p_x p_y)$ . This is convex in  $p$  and  $q$ .

If  $\mathbb{E}[\frac{\partial}{\partial x} \ln p_{y;x}] = 0$  and  $\mathbb{E}[\frac{\partial^3}{\partial x^3} \ln p_{y;x}] \neq \pm\infty$ , then  $D(p_y(\cdot|x)||p_y(\cdot|x+\delta)) = \frac{\log(e)}{2} J_y(x) \delta^2 + o(\delta^2)$ .

## Information geometry (14)

**I-projection**  $p^* = \operatorname{argmin}_{p \in \mathcal{P}} D(p||q)$ . If  $p^*$  is projection of  $q$  onto closed/convex  $\mathcal{P}$ , then  $D(p||q) \geq D(p||p^*) + D(p^*||q)$  for any  $p \in \mathcal{P}$ .

**Linear family:** family s.t.  $\{p : \mathbb{E}_p[t_i(y)] = \bar{t}_i\}$ . For  $p_1, p_2 \in \mathcal{P}$ ,  $\lambda p_1 + (1-\lambda)p_2 \in \mathcal{P}$  for any  $\lambda \in \mathbb{R}$ . Equality holds in Pythagorean Theorem if  $\mathcal{P}$  is linear.

**Orthogonal family:** given  $p^*$  in linear  $\mathcal{P}$  with parameters  $t(\cdot)$ ,  $\mathcal{E}_t(p^*)$  is all  $p$  projecting to  $p^*$ .

$$\mathcal{E}_t(p^*) = \{q \mid q(y) = p^*(y) \exp(\mathbf{x}^\top t(y) - \alpha(\mathbf{x}))\}.$$

## Modeling as inference (15)

**Mixture models**  $q_w(y) = \sum_{x \in \mathcal{X}} w(x) p_y(y|x)$ ;  $w$  is a “prior”. **Optimality:** for any distribution  $q(\cdot)$ , there exist weights  $w$  such that  $q_w$  is strictly better:  $D(p_y(\cdot|x)||q_w) \leq D(p_y(\cdot|x)||q)$  for all  $x$ .

Minimax perspective:  $\min_{q(\cdot)} \max_x D(p_y(\cdot|x)||q)$ .

**Redundancy-capacity:**  $\max \min = \min \max$ . Note the inner maximization is

$$\max_x D(p_y(\cdot|x)||q) = \max_w \sum_x w(x) D(p_y(\cdot|x)||q).$$

Interpret  $w$  as a **least informative prior**  $p_x^* = \operatorname{argmax}_{p_x} I(x; y)$ . **Model capacity** is  $C = \max_{p_x} I(x; y) \leq \log |\mathcal{X}|$ . **Equidistance property:** at mixture  $q^*$  and weights  $w^*$ ,  $D(p_y(\cdot|x)||q^*) \leq C \forall x$ , equality when  $w^*(x) > 0$ .

## Continuous information theory (16)

- Entropy  $h(x) = -\int p(x) \log p(x) dx$ . Conditional  $h(x|y) = -\int p_y(y) h(x|y=y) dy$ .
- Mutual information  $I(x; y) = h(x) - h(x|y)$ .
- Divergence  $D(p||q) = \int p(x) \log \frac{p(x)}{q(x)} dx$ .

For Gaussians,  $h(x) = \frac{\log e}{2} (\ln |2\pi\Lambda| + \operatorname{tr}(I))$ ,  $h(x|y) = \frac{1}{2} \log |2\pi e(\Lambda_x - \Lambda_{xy}\Lambda_y^{-1}\Lambda_{xy}^\top)|$ .

## Some useful theorems (11 & misc.)

**Jensen's inequality:** For concave  $\varphi(\cdot)$  and random variable  $v$ ,

$$\mathbb{E}[\varphi(v)] \geq \varphi(\mathbb{E}[v]).$$

**Csiszár's inequality:** for convex  $f$ ,

$$\sum_{i=1}^N b_i f\left(\frac{a_i}{b_i}\right) \geq \left(\sum_{i=1}^N b_i\right) f\left(\frac{\sum_{i=1}^N a_i}{\sum_{i=1}^N b_i}\right).$$

**Log-sum:**  $\sum a_i \log \frac{a_i}{b_i} \geq (\sum a_i) \log \frac{\sum a_i}{\sum b_i}$ .

**Gibbs' inequality:** let r.v.  $v \sim p(\cdot)$ . For any distribution  $q(\cdot)$ ,  $\mathbb{E}_p[\log p(v)] \geq \mathbb{E}_p[\log q(v)]$ .

**Lagrange multipliers:**  $\nabla f(\hat{z}) + \lambda^\top \nabla g(\hat{z}) = 0$ .

**Taylor's Theorem:** for some  $t'$  between  $t$  and  $t_0$ ,

$$f(t) = \sum_{j=0}^J (t-t_0)^j \frac{f^{(j)}(t_0)}{j!} + (t-t_0)^{J+1} \frac{f^{(J+1)}(t')}{(J+1)!}$$