

Two class problem

$y \in \{0, 1\}$ disease, $x \in \{0, 1\}$ diagnosis

sample $(x_i, y_i) \stackrel{\text{i.i.d.}}{\sim} P(x, y)$, $i = 1 \dots n$
 \uparrow population frequency

$$P(x, y) = \frac{1}{n} \sum_{i=1}^n 1(x_i = x) 1(y_i = y)$$

\uparrow sample frequency

Two "models": $P(x | y=1)$ vs. $P(x | y=0)$

"Bayes factor": $\frac{P(x | y=1)}{P(x | y=0)} \stackrel{\text{e.g.}}{=} 3.0 \Rightarrow$ "x is 3x more likely under $y=1$ "
i.e. likelihood ratio (LR)

If $LR > 1$, do we conclude $y=1$?

What about the "base rate" $P_r(y=1)$?

"Prior" model: $P(y=1)$

The correct probabilistic query:

$$P(y=1 | x) = \frac{P(y=1 | x) P(y=1)}{P(x)}$$

\uparrow "posterior" \uparrow "likelihood" \uparrow "prior" \uparrow "marginal likelihood"

Bayes rule

Posterior odds

LR

Prior odds

$$\frac{P(y=1 | x)}{P(y=0 | x)} = \frac{P(x | y=1)}{P(x | y=0)} \times \frac{P(y=1)}{P(y=0)}$$

e.g. $P(x=1 | y=1) = 1.0$ true positive rate of diagnostic

$P(x=1 | y=0) = 0.1$ false positive rate of diagnostic

$P(y=1) = 0.001$ rare disease

$$\hookrightarrow \text{posterior odds} = \frac{1.0}{0.1} \times \frac{0.001}{0.999} \approx 0.01$$

"data x is 100-to-1 against disease"

Moral: Priors matter! See "Lady drinking tea" (Berger, 1985)

Decision theory: Say posterior odds are > 1 (e.g. 10).

Do we treat the patient?

decision rule: $\hat{y}(x) \in \{0, 1\}$

loss function: $\ell(y, \hat{y}(x))$

For binary classification:

$\hat{y}=0$	$\hat{y}=1$
$\ell(0,0)$ TN	$\ell(0,1)$ FP $y=0$
$\ell(1,0)$ FN	$\ell(1,1)$ TP $y=1$

l_{TN} and l_{FP} should be $\leq l_{FN}$ and l_{FP}

but $l_{FN} \neq l_{FP}$ (not necessarily)

$$\text{e.g. } \hat{y}(x) = \begin{cases} 1 & \text{if } P(y=1|x) \geq 0.5 \\ 0 & \text{otherwise} \end{cases}$$

Integrated risk:

$$\begin{aligned} r(\hat{y}(\cdot)) &= \mathbb{E}_{x,y} [l(\hat{y}(x), y)] \\ &= \mathbb{E}_{(x,y) \sim P_r} [\dots] \\ &\quad \leftarrow \text{population average} \end{aligned}$$

Optimal rule

$$\begin{aligned} \hat{y}^*(\cdot) &= \underset{\hat{y}(\cdot)}{\operatorname{argmin}} r(\hat{y}(\cdot)) \\ &= \underset{\hat{y}(\cdot)}{\operatorname{argmin}} \mathbb{E}_x \mathbb{E}_{y|x} [l(\hat{y}(x), x)] \end{aligned}$$

The optimal rule "always does the right thing":

$$\hat{y}^*(x) = \underset{a}{\operatorname{argmin}} \mathbb{E}_{y|x} [l(a, x)]$$

\leftarrow "posterior"

In other words we "only" need to know $P_r(y|x)$

$\hookrightarrow \hat{y}^*(x)$ is the "Bayes decision rule"

$\hookrightarrow r(\hat{y}^*(\cdot))$ is the "Bayes risk"

Frequentist justification for "Bayesian inference" = posterior computation

Modeling: Discriminative vs. generative

We almost never know $\Pr(y|x) \dots$

... especially true for higher-dimensional $x \in \mathbb{R}^p$

"making assumptions about" $\Pr(y|x)$ vs. $\Pr(x|y) \Pr(y)$

$X_1 \dots X_p$ are symptoms, $X_i \in \{0, 1\}$ e.g. "is coughing"

$X \equiv X_{1:p}$ takes 2^p possible values

e.g. $p=40 \Rightarrow 2^{40} \approx 1$ trillion

no hope to estimate $\Pr(y|x)$ $\forall x$ without assumptions

Discriminative example

Assume: $\Pr(y|x) = \frac{1}{1 + \exp(-\beta^T x)}$ for some $\beta \in \mathbb{R}^p$
(logistic regression)

Allows us to "share information" across similar x .

Generative example: Naive Bayes

Assume: $\Pr(x_{1:p} | y) = \prod_{j=1}^p \Pr(x_j | y)$

i.e. $x_j \perp x_{j'} | y$ conditionally indep.

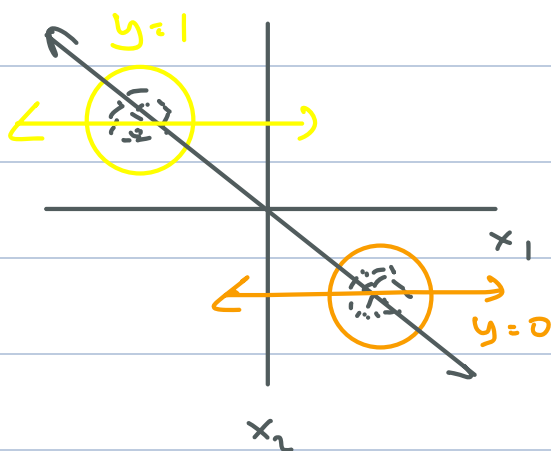
$$\hookrightarrow P_r(y|x) \stackrel{\text{Bayes}}{=} \frac{P(y) P(x|y)}{P_r(x)} \stackrel{\text{naive}}{=} \frac{P_r(y) \prod_j P(x_j|y)}{P_r(x)}$$

What about $P(x) = P(x_1, \dots, x_p)$? $\neq \prod_j P(x_j)$

$$= \sum_y P(y) P(x|y) \quad = \sum_y P(y) \prod_j P(x_j|y)$$

\sum_y marginalization

Note: $x_j \perp\!\!\!\perp x_{j'} \mid y \not\Rightarrow x_j \perp\!\!\!\perp x_{j'}$



related: Simpson's paradox

Takeaway: Only need to estimate $P(y)$ and $P(x_j|y)$.

Again, "information sharing". \uparrow $2p$ parameters vs. 2^p

This model is (often) "wrong" but (often) "useful". See Box.

Data sparsity:

$$P(x_j=1 \mid y=1) = \theta_{j11} \leftarrow \frac{\sum_i 1(x_{i1}=1, y_i=1)}{\sum_i 1(y_i=1)}$$

What if $\theta_{j11} = 1$ (or 0)?

e.g. $P(y^{n+1}=1 \mid x_1^{n+1}=1, x_2^{n+1}=1, \dots)$

$$= \frac{P(y=1) \prod_i \theta_{i,1}}{P(y=1) \prod_i \theta_{i,1} + P(y=0) \prod_i \theta_{i,0}}$$

$$\text{say } \theta_{1,1} = 1 \quad \theta_{1,0} = 0 \quad \theta_{2,1} = 0 \quad \theta_{2,0} = 1$$

$$L = \frac{0}{0} \quad \ddots$$

Laplace smoothing

$$\theta_{i,1} \leftarrow \frac{\#(x_i=1, y=1) + \alpha}{\#(y=1) + 2\alpha}$$

Pierre-Simon Laplace (1700s) proposed this for the "Sunrise problem"

see also: Bertrand Russell's "problem of induction" (Bertrand's chicken)

Justification:

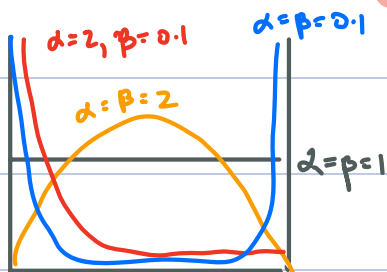
$$x^i \stackrel{iid}{\sim} \text{Bernoulli}(\theta) \quad i=1 \dots n$$

$$\hat{\theta}_{MLE} = \underset{\theta}{\text{argmax}} \prod_i P(x^i | \theta) = \bar{x}$$

Prior: $\theta \sim \text{Beta}(\alpha, \beta)$

$$P(\theta | \alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \mathbb{1}(\theta \in (0,1))$$

"normalizing constant"
"kernel"



Alternative parametrization:

$$\mu = E[\theta | \alpha, \beta] = \frac{\alpha}{\alpha + \beta}$$

$$\kappa = \alpha + \beta \leftarrow \text{"concentration"}$$

Posterior:

$$P(\theta | x_{1:n}) = \frac{P(\theta) P(x_{1:n} | \theta)}{P(x_{1:n})} \longrightarrow = \int d\theta P(\theta) P(x | \theta)$$

$$= \int \frac{\text{Beta}(\theta; \alpha, \beta) \prod_i \text{Bern}(x_i; \theta)}{d\theta}$$

$P(x_{1:n})$ is the "normalizing constant" of the posterior

$$\propto \text{Beta}(\theta; \alpha, \beta) \prod_i \text{Bern}(x_i; \theta)$$

$$= \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \prod_i \theta^{x_i} (1-\theta)^{1-x_i}$$

$$\propto \theta^{\alpha + \sum_i x_i - 1} (1-\theta)^{\beta + n - \sum_i x_i - 1}$$

This is the kernel of another Beta distribution.

$$\hookrightarrow P(\theta | x_{1:n}) = \text{Beta}(\alpha + \sum_i x_i, \beta + n - \sum_i x_i)$$

No need to "solve" for the normalizing constant explicitly.

Why did this happen? Beta-Bernoulli conjugacy.

Beta is the conjugate prior for the Bernoulli.

Bayesian updating

$$\text{Beta}(\underbrace{\alpha + \sum_i x_i}_{\# \text{ 1s}}, \underbrace{\beta + n - \sum_i x_i}_{\# \text{ 0s}})$$

α, β are "pseudocounts"

$$\mathbb{E}[\theta | x_{1:n}, \alpha, \beta] = \frac{\alpha + \sum_i x_i}{\alpha + \beta + n} \quad \leftarrow \text{Laplace smoothing estimator!}$$

$$= \frac{k}{n} \cdot \frac{\alpha}{\alpha + n} + \frac{\sum x_i}{\alpha + n} \cdot \frac{n}{n} = \left(\frac{\alpha}{\alpha + n}\right) \mu + \left(\frac{n}{\alpha + n}\right) \hat{\theta}_{\text{MLE}}$$

$\alpha + \beta$ is the strength of the prior / inductive bias.

Subjectivist interpretation of probability

$P(\theta) = \text{Beta}(\theta; \alpha, \beta)$ is not a "frequency"
... it is a "degree of belief".

e.g. $P(\text{Bulls win on Thursday})$
 $P(\text{sun comes up tomorrow})$
 $P(\text{Riemann's conjecture is True})$

The Bayesian vs frequentist debate was over the interpretation (among other things).