

Safe Probability



Peter Grünwald



Centrum Wiskunde & Informatica – Amsterdam
Mathematisch Instituut – Universiteit Leiden

Prelude: Kelly Gambling

- At each time i , we can buy a ticket $T_{i,1}$ that pays off \$2 iff $X_i = 1$, and a ticket $T_{i,0}$ that pays off \$2 iff $X_i = 0$. Both tickets cost \$1
- A **gambling strategy** in this game is a function $q : \bigcup_{n \geq 0} \{0,1\}^n \rightarrow \Delta_2$ and thus defines a probability distr. on $\{0,1\}^\infty$ via setting

$$\tilde{P}(X_i = 1 \mid X^{i-1}) := q(1 \mid X^{i-1}) \quad ; \quad \tilde{P}(X^n) := \prod_{i=1}^n \tilde{P}(X_i \mid X^{i-1})$$

- If we follow such a strategy and start with \$1, our capital after n rounds will be $\tilde{P}(X^n) \cdot 2^n$

How to design a gambling strategy?

- A **gambling strategy** in this game is formally equivalent to a probability distribution \tilde{P} on infinite sequences. **Which strategy should we adopt?**
- Strict Subjective Bayesian**: think very long about the situation, come up with a subjective distribution P^* , and then play the distribution \tilde{P} maximizing expected gain (we may have $\tilde{P} \neq P^*$)
- Imprecise Probabilist**: come up with a set of distributions \mathcal{P}^* , and then play the distribution \tilde{P} optimal relative to \mathcal{P}^* , with optimality defined relative to some additional criterion (which one?)

Prelude: Kelly Gambling

- Suppose we observe sequence X_1, X_2, \dots of 0s and 1s
- At each point in time i , we can buy a ticket $T_{i,1}$ that pays off \$2 iff $X_i = 1$, and a ticket $T_{i,0}$ that pays off \$2 iff $X_i = 0$. Both tickets cost \$1
- Crucially: we are allowed to **divide** our capital any way we like and **re-invest** our capital at each point in time
 - e.g. By putting 50% of your capital at time i on $T_{i,1}$ and 50% on $T_{i,0}$ you make sure that your capital remains the same

How to design a gambling strategy?

- A **gambling strategy** in this game is formally equivalent to a probability distribution \tilde{P} on infinite sequences. **Which strategy should we adopt?**

How to design a gambling strategy?

- Strict Subjective Bayesian**: determine subjective P^* , and then play optimal \tilde{P} (we may have $\tilde{P} \neq P^*$)
- Imprecise**: determine set \mathcal{P}^* and play "optimal" \tilde{P}
- Information Theorist**: pick any gambling strategy which you think might gain you a lot. E.g. if you think frequency might converge to $p \neq 0.5$, you might play **Laplace rule of succession...**

$$\tilde{P}(X_{n+1} = 1 \mid X^n) = \frac{n+1}{n+2} \quad \tilde{P}(X^n) = \int_0^1 p^{n+1} (1-p)^{n_0} dp$$

How to design a gambling strategy?

- **Strict Subjective Bayesian**: determine subjective P^* , and then play optimal \tilde{P} (we may have $\tilde{P} \neq P^*$)
- **Imprecise**: determine set \mathcal{P}^* and play "optimal" \tilde{P}
- **Information Theorist**: pick **any gambling strategy** which you think might gain you a lot. E.g. if you think frequency might converge to $p \neq 0.5$, you might play **Laplace rule of succession...**

$$\tilde{P}(X_{n+1} = 1 | X^n) = \frac{n_1+1}{n+2} \quad \tilde{P}(X^n) = \int_0^1 p^{n_1} (1-p)^{n_0} dp$$

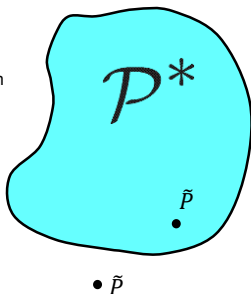
...if your hypothesis about frequency is correct, you gain **exponential amount of money** even if at the same time you think data are **not Bernoulli (or not even stationary)**

Starting Point

- Adopting a Bayesian predictive distribution like the Laplace Rule of Succession if you think data are not Bernoulli is o.k. (**and I think, rational!**) for some prediction tasks...
 - Sequential gambling, **Data Compression**
 ...but not for others:
 - 0/1-loss prediction (no fractional bets!) when you are only asked to predict X_i in the situation that $X_{i-1} = 1$
- I want to design a theory which can cope with such 'partially useable' distributions

A Middle Ground between strict Bayes and imprecise probability

- Set of distrs \mathcal{P}^* has unique **representative**, as in 'objective Bayes', fiducial inference, Maximum Entropy, data compression...
- One absolutely crucial difference: we **restrict** use of \tilde{P} to subset of all possible prediction tasks: we know in advance that \tilde{P} **should not be taken to seriously**
- **Provides unifying and demistifying view**



Menu

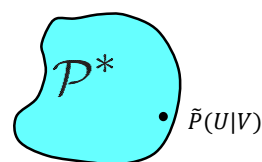
1. The Setting
2. Definition 1, Example 1: **Dilation**
3. Definition 2, Example 1 cont.
4. Definition 3-4, Example 2: **Calibration**
5. Example 3: **Fiducial** Distributions
6. Desert: **Monty Hall** Problem, Decision Safety

The Setting

- Let \mathcal{P}^* be a **set of distributions** on a space Ω , representing Decision-Maker (DM)'s uncertainty about a domain
- DM has to make predictions/assertions about some U (or a function thereof), upon observing V . Both U and V are RVs (random variables) on Ω , taking values in \mathcal{U} and \mathcal{V} , resp.
- She does so using a **pragmatic distribution** $\tilde{P}(U|V)$, defined as a **conditional distribution** of U given V , i.e. a function mapping each $v \in \mathcal{V}$ to a distribution $\tilde{P}(U | V = v)$ on \mathcal{U}
 - Whenever \mathcal{U} finite, we think of $\tilde{P}(U) = (\tilde{p}(u_1), \dots, \tilde{p}(u_m))^T$ as a column vector

The Setting

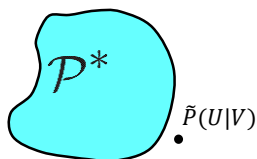
- A **Bayesian** would have $\mathcal{P}^* = \{P^*\}$ a singleton and could then set $\tilde{P}(U | V) = P^*(U | V)$
- Note that \mathcal{P}^* is a distribution on Ω , inducing a joint $P^*(U, V)$ which in turn induces $P^*(U | V)$, while $\tilde{P}(U | V)$ is directly defined as a conditional (hence \tilde{P} in picture to be taken with grain of salt)



The Setting

- A **Bayesian** would have $\mathcal{P}^* = \{P^*\}$ a singleton and could then set $\tilde{P}(U|V) = P^*(U|V)$
- We have to do something else – sometimes eqv. to conditioning on a special element of \mathcal{P}^* , sometimes really different...

\tilde{P} is really a probability update rule!!



First Definition: Weak Safety

- We say that $\tilde{P}(U|V)$ is **safe** for $\langle U \rangle | \langle V \rangle$ if for all $P^* \in \mathcal{P}$:

$$E_{P^*}[U] = E_{V \sim P^*} E_{U \sim \tilde{P}|V}[U]$$

- i.e.

$$\sum_{u \in \mathcal{U}} p^*(u) \cdot u = \sum_{v \in \mathcal{V}} p^*(v) \sum_{u \in \mathcal{U}} \tilde{p}(u|v) \cdot u$$

First Definition

- We say that $\tilde{P}(U|V)$ is **safe** for $\langle U \rangle | \langle V \rangle$ if for all $P^* \in \mathcal{P}$:

$$E_{P^*}[U] = E_{V \sim P^*} E_{U \sim \tilde{P}|V}[U]$$

- i.e.

$$E_{V \sim P^*} E_{U \sim P^*|V}[U] = E_{V \sim P^*} E_{U \sim \tilde{P}|V}[U]$$

First Definition

- We say that $\tilde{P}(U|V)$ is **safe** for $\langle U \rangle | \langle V \rangle$ if for all $P^* \in \mathcal{P}$:

$$E_{P^*}[U] = E_{V \sim P^*} E_{U \sim \tilde{P}|V}[U]$$

- i.e. we can expect our expectation of U to be ‘correct’ (in a relative sense)
- we will usually want somewhat stronger versions of ‘safety’

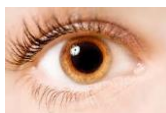
First Example: Dilation

$$\mathcal{U} = \mathcal{V} = \{0, 1\}, \mathcal{P} = \{P^* : E_{P^*}[U] = 0.9\}$$

- **Given:** marginal probability of U . U may depend on V , but we have no idea how
- **Task:** predict U given V .
- Suppose we observe $V = 0$. Now conditional probability could be anything...

$$\mathcal{P}(U|V=0) := \{P|V=0 : P \in \mathcal{P}\} = [0, 1]$$
- Similarly if we observe $V = 1$:

$$\mathcal{P}(U|V=1) := \{P|V=1 : P \in \mathcal{P}\} = [0, 1]$$



Dilation

Seidenfeld & Wasserman, '93

Before observing V we had precise probability $P^*(U = 1)$
after we only know $P^*(U = 1 | V = v)$ is in large superset

“extra information \Rightarrow less knowledge
 no matter what you observe!”

Ignoring instead of Dilating

- Pointwise conditioning gives dilation
- Instead we may decide to **ignore** V , i.e. act as if U and V are **independent**, and predict with the **pragmatic** distribution

$$\tilde{P}(U = 1 | V = 1) = \tilde{P}(U = 1 | V = 0) = 0.9$$

- Proposition:** $\tilde{P}(U|V)$ is **safe** for $\langle U \rangle | \langle V \rangle$

$$E_{P^*}[U] = E_{V \sim P^*} E_{U \sim \tilde{P}|V}[U]$$

- i.e. $P^*(U = 1) = E_{V \sim P^*}[\tilde{P}(U = 1 | V)]$

First Example of ‘Safety’

- REALITY:** U may be **dependent** on V
- PRAGMATICS:** we nevertheless decide to predict U with a distribution that assumes U and V are **independent**
- Our predictions will be **just as accurate as we would expect them to be if our pragmatic distribution \tilde{P} were ‘correct’**
- ...as long as we only use \tilde{P} only for certain, not all prediction tasks...

Definition 2, Preparation

- We write $X \rightsquigarrow Y$ if there exists a function ϕ such that $\phi(X) \equiv Y$ (“ X **determines** Y ”)
- $\tilde{P}(U|V)$ can be used to predict not just U , but also any U' determined by (U, V) , i.e. with $(U, V) \rightsquigarrow U'$:

$$\tilde{P}(U' = u' | V = v) := \sum_{u \in \mathcal{U}: \phi(u, v) = u'} \tilde{P}(U = u | V = v)$$

Definition 2

- Recall: $\tilde{P}(U|V)$ is **safe** for $\langle U' \rangle | \langle V \rangle$ if $(U, V) \rightsquigarrow U'$ and for all $P^* \in \mathcal{P}$:

$$E_{P^*}[U'] = E_{V \sim P^*} E_{U' \sim \tilde{P}|V}[U']$$

- We say that $\tilde{P}(U|V)$ is **safe** for $U | \langle V \rangle$ if **for all** U' with $(U, V) \rightsquigarrow U'$, all $P^* \in \mathcal{P}$:

$$E_{P^*}[U'] = E_{V \sim P^*} E_{U' \sim \tilde{P}|V}[U']$$

Definition 2

- Recall: $\tilde{P}(U|V)$ is **safe** for $\langle U' \rangle | \langle V \rangle$ if $(U, V) \rightsquigarrow U'$ and for all $P^* \in \mathcal{P}$:

$$E_{P^*}[U'] = E_{V \sim P^*} E_{U' \sim \tilde{P}|V}[U']$$

- We say that $\tilde{P}(U|V)$ is **safe** for $U | \langle V \rangle$ if **for all** U' with $(U, V) \rightsquigarrow U'$, all $P^* \in \mathcal{P}$:

$$E_{P^*}[U'] = E_{V \sim P^*} E_{U' \sim \tilde{P}|V}[U']$$

$$P^*(U) = E_{V \sim P^*}[\tilde{P}(U | V)]$$

Example 1(b) - dilation again

$$\mathcal{U} = \mathcal{V} = \{0, 1, 2\}, \mathcal{P} = \{P^* : E_{P^*}[U] = 0.9\}$$

- Task:** predict U given V .
- Again we decide to ignore V and set e.g. for all $v \in \mathcal{V}$:

$$\tilde{P}(U = 1 | V = v) = 0.9, \tilde{P}(U = 2 | V = v) = 0$$

- Then \tilde{P} is **safe** for $\langle U \rangle | \langle V \rangle$ but not for $U | \langle V \rangle$

Example 1(c)

$$\mathcal{U} = \mathcal{V} = \{0, 1, 2\}, \mathcal{P} = \{P^* : P^*(U) = (0, 0.3, 0.3)^T\}$$

- **Task:** predict U given V .
- Again we decide to ignore V and set e.g. for all $v \in \mathcal{V}$:

$$\tilde{P}(U = 1 | V = v) = 0.9, \tilde{P}(U = 2 | V = v) = 0$$

- Then, **again**, \tilde{P} is safe for $\langle U \rangle | \langle V \rangle$ but not for $U | \langle V \rangle$

Example 1(c): use the marginal

$$\mathcal{U} = \mathcal{V} = \{0, 1, 2\}, \mathcal{P} = \{P^* : P^*(U) = (0, 0.3, 0.3)^T\}$$

- **Task:** predict U given V .
- Again we decide to ignore V and set e.g. for all $v \in \mathcal{V}$:

$$\tilde{P}(U | V = v) = P^*(U)$$

- Then \tilde{P} is safe for $\langle U \rangle | \langle V \rangle$ and **also** for $U | \langle V \rangle$

Definition 3, Preparation

- Recall: $\tilde{P}(U|V)$ is **safe** for $\langle U' \rangle | \langle V \rangle$ if $(U, V) \rightsquigarrow U'$ and for all $P^* \in \mathcal{P}$:

$$E_{P^*}[U'] = E_{V \sim P^*} E_{U' \sim \tilde{P}|V}[U']$$

- Leave out ' $(U, V) \rightsquigarrow U'$ ' part from now on, for brevity

Definition 3

- Recall: $\tilde{P}(U|V)$ is **safe** for $\langle U' \rangle | \langle V \rangle$ if for all $P^* \in \mathcal{P}$:

$$E_{P^*}[U'] = E_{V \sim P^*} E_{U' \sim \tilde{P}|V}[U']$$

- We say that $\tilde{P}(U|V)$ is **safe** for $\langle U' \rangle | V$ if for all $P^* \in \mathcal{P}$:

$$E_{P^*}[U' | V] \equiv E_{\tilde{P}}[U' | V]$$

Definition 3

- Recall: $\tilde{P}(U|V)$ is **safe** for $\langle U' \rangle | \langle V \rangle$ if for all $P^* \in \mathcal{P}$:

$$E_{P^*}[U'] = E_{V \sim P^*} E_{U' \sim \tilde{P}|V}[U']$$

- We say that $\tilde{P}(U|V)$ is **safe** for $\langle U' \rangle | V$ if for all $P^* \in \mathcal{P}$:

$$E_{P^*}[U' | V] = E_{\tilde{P}}[U' | V]$$

- Our expectation of U' is (relatively) correct

Definition 3, 3b

- Recall: $\tilde{P}(U|V)$ is **safe** for $\langle U' \rangle | \langle V \rangle$ if for all $P^* \in \mathcal{P}$:

$$E_{P^*}[U'] = E_{V \sim P^*} E_{U' \sim \tilde{P}|V}[U']$$

- We say that $\tilde{P}(U|V)$ is **safe** for $\langle U' \rangle | V$ if for all $P^* \in \mathcal{P}$:

$$E_{P^*}[U' | V] = E_{\tilde{P}}[U' | V]$$

- We say that $\tilde{P}(U|V)$ is **safe** for $U' | V$ if for all $P^* \in \mathcal{P}$:

$$P^*[U' | V] \equiv \tilde{P}[U' | V]$$

i.e. P^* is unique and \tilde{P} is almost surely 'correct'

Definition 3c...

- We can now also combine definitions, e.g. $\tilde{P}(U|V)$ is **safe** for $\langle U' \rangle | \langle V \rangle, W$ if(details omitted)...

Ex. 2, Calibration: preparation

- Recall: $\tilde{P}(U|V)$ can be used to predict not just U , but also any U' determined by (U, V) , i.e. with $(U, V) \rightsquigarrow U'$

$$\tilde{P}(U' = u' | V = v) := \sum_{(u,v) \in \mathcal{U} \times \mathcal{V} : \phi(u,v)=u'} \tilde{P}(U = u | v)$$

$$\tilde{P}(U' = u' | V = v_1) = \tilde{P}(U' = u' | V = v_2)$$

Ex. 2, Calibration: preparation

- Recall: $\tilde{P}(U|V)$ can be used to predict not just U , but also any U' determined by (U, V) , i.e. with $(U, V) \rightsquigarrow U'$
- Similarly, $\tilde{P}(U|V)$ can also be used to predict not just given V , but also given any V' with $V \rightsquigarrow V'$ and **extra condition** that for all $v_1, v_2, v_1 \neq v_2$:

$$\phi_{V \rightsquigarrow V'}(v_1) = \phi_{V \rightsquigarrow V'}(v_2) \Rightarrow \tilde{P}(U | V = v_1) = \tilde{P}(U | V = v_2)$$

- For such V' , $\tilde{P}(U|V')$ is well-defined
- Example:** earlier \tilde{P} that treated U, V as independent:
 $V \in \{0, 1\}, V' \equiv 0$

Ex. 2, preparation

- $\tilde{P}(U|V)$ can be used to predict not just given V , but also given any V' with $V \rightsquigarrow V'$ and **extra condition** that for all $v_1, v_2, v_1 \neq v_2$:

$$\begin{aligned} \phi_{V \rightsquigarrow V'}(v_1) &= \phi_{V \rightsquigarrow V'}(v_2) \\ \Rightarrow \tilde{p}(U | V = v_1) &= \tilde{P}(U | V = v_2) \end{aligned}$$

- compact restatement:** $\tilde{P}(U|V)$ can also be used to predict (i.e. induces a unique definition of $\tilde{P}(U|V')$) based on any V' with $V \rightsquigarrow V'$

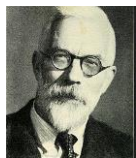
Ex. 2., Calibration

- We say that $\tilde{P}(U|V)$ is **strongly calibrated** for $U' | V$ if it is **safe** for $U' | \tilde{P}(U|V)$!
...i.e. for all $P^* \in \mathcal{P}, \vec{p} \in \text{Range}(P(U | V))$
 $P^*(U' | \tilde{P}(U | V) = \vec{p}) = \vec{p}$



Ex. 2., Calibration

- We say that $\tilde{P}(U|V)$ is **strongly calibrated** for $U' | V$ if it is **safe** for $U' | \tilde{P}(U|V)$!
...i.e. for all $P^* \in \mathcal{P}, \vec{p} \in \text{Range}(P(U | V))$
 $P^*(U' | \tilde{P}(U | V) = \vec{p}) = \vec{p}$
- Ex.: a **weather forecaster** predicts daily precipitation probabilities $\tilde{P}(U|V)$, based on measurements of air pressure and temperature taken all over the world
– so V is a giant vector. WF will probably not be able to give accurate predictions given the air pressure in Honolulu, although his predictions do depend thereon. We don't mind this, but we do want her to be calibrated!



Ex. 3, Fiducial Distributions

- Determining a distribution on parameters $\tilde{p}(\theta | X^n)$ without a prior, i.e. **cooking a Bayesian omelet without breaking the Bayesian eggs**
- Introduced by **Fisher** (1935, *Annals of Eugenics*)
Once almost on a par with Bayes and frequentist approaches; but turned out to suffer severe difficulties – e.g., **Seidenfeld '92**
- Yet it is making a small comeback under the name **confidence distributions** (Hjort & Schweder, 2000)

Fiducial Distributions

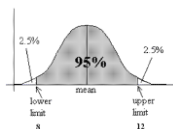
- Simple Example: normal location family
- $$\mathcal{M} = \{p_\theta | \theta \in \mathbb{R}\}, p_\theta(X^n) = \prod_{i=1}^n p_\theta(X_i), p_\theta(X_i) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(X_i - \theta)^2}$$
- Fisher observed that the density of the ML estimator $\hat{\theta} = n^{-1} \sum_{i=1}^n X_i$ satisfies
- $$p_\theta(\hat{\theta}) = \sqrt{\frac{n}{2\pi}} \cdot e^{-\frac{\sqrt{n}}{2} \cdot (\hat{\theta} - \theta)^2}$$
- which is **symmetric** in $\theta, \hat{\theta}$... so that for each $X^n \in \mathbb{R}^n$ $\tilde{p}(\theta | X^n) := p_\theta(\hat{\theta})$ must give a distribution on θ ...

Fiducial Distributions

- Simple Example: normal location family
- $$\mathcal{M} = \{p_\theta | \theta \in \mathbb{R}\}, p_\theta(X^n) = \prod_{i=1}^n p_\theta(X_i), p_\theta(X_i) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(X_i - \theta)^2}$$
- Fisher observed that the density of the ML estimator $\hat{\theta} = n^{-1} \sum_{i=1}^n X_i$ satisfies
- $$p_\theta(\hat{\theta}) = \sqrt{\frac{n}{2\pi}} \cdot e^{-\frac{\sqrt{n}}{2} \cdot (\hat{\theta} - \theta)^2}$$
- which is **symmetric** in $\theta, \hat{\theta}$... so that for each $X^n \in \mathbb{R}^n$ $\tilde{p}(\theta | X^n) := p_\theta(\hat{\theta})$ must give a distribution on θ ...
- ...Fisher now boldly treated this is a sort-of posterior...

Fiducial Distribution

- Can do similar reversal for other 1-parameter distributions.
 - For scale and location families, the fiducial distr is equal to the Bayes' posterior with the improper Jeffreys' prior
 - For other families, no 100% Bayes interpretation (Lindley, Seidenfeld)
 - For **Bayesians this seems flawed**: there must be a prior
 - For **Frequentists this seems flawed**: θ is fixed, not a random variable!!
- $$\tilde{p}(\theta | X^n) := \sqrt{\frac{n}{2\pi}} \cdot e^{-\frac{\sqrt{n}}{2} \cdot (\hat{\theta} - \theta)^2}$$



Fiducial Distributions and Confidence

- It has long been known that **fiducial distributions are "o.k." if used to determine confidence intervals...**
suppose X_1, X_2, \dots i.i.d. $\sim P_{\theta^*}$, for any $\theta^* \in \mathbb{R}$
Set $\theta^+ = \theta^+(X^n)$ and $\theta^- = \theta^-(X^n)$ such that
- $$\tilde{P}(\theta \leq \theta^-(X^n) | X^n) = 0.025$$
- $$\tilde{P}(\theta \geq \theta^+(X^n) | X^n) = 0.025$$
- Then: $P_{\theta^*}(\theta^* \notin [\theta^-, \theta^+]) = 0.05$

Fiducial Distributions and Confidence

- It has long been known that **fiducial distributions are "o.k." if used to determine confidence intervals...**
- ...but are not o.k. "in general" (but what exactly does this mean? And what are they o.k. for?)

Using Fiducial Distributions **Safely**

- Let \mathcal{P} be **any** set of distributions on $\Theta \times \mathbb{R}^n = \mathcal{U} \times \mathcal{V}$ such that for all $\theta^* \in \mathbb{R}, x^n \in \mathbb{R}^n : p(x^n|\theta^*) := p_{\theta^*}(x^n)$ (\mathcal{P} may e.g. only contain degenerate distr on Θ)
- We define some 'pragmatic posterior' $\tilde{P}(\theta | X^n)$ and...
- DEF:** we say that $\tilde{P}(\theta | X^n)$ is **fiducially safe** if it is safe for $\tilde{F}(\cdot | X^n) | \langle X^n \rangle$, where

$$\tilde{F}(u | X^n) = \tilde{P}(\theta \leq u | X^n)$$

is the distribution function of $\tilde{P}(\theta | X^n)$

Using Fiducial Distributions **Safely**

- Let \mathcal{P} be **any** set of distributions on $\Theta \times \mathbb{R}^n = \mathcal{U} \times \mathcal{V}$ such that for all $\theta^* \in \mathbb{R}, x^n \in \mathbb{R}^n : p(x^n|\theta^*) := p_{\theta^*}(x^n)$
- DEF:** we say that $\tilde{P}(\theta | X^n)$ if it is safe for $\tilde{F}(\cdot | X^n) | \langle X^n \rangle$
- PROP:** if defined in the usual way, 'fiducial' distributions are fiducially safe

Using Fiducial Distributions **Safely**

- Let \mathcal{P} be **any** set of distributions on $\Theta \times \mathbb{R}^n = \mathcal{U} \times \mathcal{V}$ such that for all $\theta^* \in \mathbb{R}, x^n \in \mathbb{R}^n : p(x^n|\theta^*) := p_{\theta^*}(x^n)$
- DEF:** we say that $\tilde{P}(\theta | X^n)$ if it is safe for $\tilde{F}(\cdot | X^n) | \langle X^n \rangle$
- PROP:** if defined in the usual way, 'fiducial' distributions are fiducially safe
- This means they can be safely used to predict any RV U' determined by $\tilde{F}(\theta | X^n)$
 - For example, $\mathbf{1}_{\theta \notin [\theta^-, \theta^+]}$ is fiducially safe to predict...

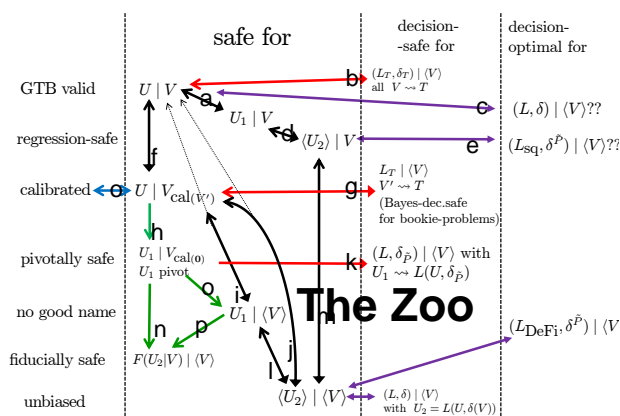
$$P_{\theta^*}(\theta^* \notin [\theta^-, \theta^+]) = E_{X^n \sim p_{\theta^*}}[\tilde{P}(\theta \notin [\theta^-, \theta^+] | X^n)] = 0.05$$

Using Fiducial Distributions **Safely**

- Let \mathcal{P} be **any** set of distributions on $\Theta \times \mathbb{R}^n = \mathcal{U} \times \mathcal{V}$ such that for all $\theta^* \in \mathbb{R}, x^n \in \mathbb{R}^n : p(x^n|\theta^*) := p_{\theta^*}(x^n)$
- DEF:** we say that $\tilde{P}(\theta | X^n)$ is **fiducially safe** for $\theta | \langle X^n \rangle$ if it is safe for $\tilde{F}(\cdot | X^n) | \langle X^n \rangle$
- PROP:** if defined in the usual way, 'fiducial' distributions are fiducially safe
- This means they can be safely used to predict any RV U' determined by $\tilde{F}(\theta | X^n)$
 - For example, $\mathbf{1}_{\theta \notin [\theta^-, \theta^+]}$ is fiducially safe to predict...
 - ...but $\mathbf{1}_{\theta \notin [12.3, 19.8]}$ is not!

Dilation-Fiducial Duality

- DILATION-REALITY:** U may be **dependent** on V
- DILATION-PRAGMATICS:** we nevertheless decide to predict U with a distribution that assumes U and V are **independent**
- FIDUCIAL-REALITY:** U may be **independent** of V, it may even be fixed – but its value is unknown
- FIDUCIAL-PRAGMATICS:** we nevertheless predict U with a distribution that assumes U and V are **dependent**



Desert: Monty Hall (3-door) Problem

Monty Hall 1970



Monty Hall



- There are three doors in the TV studio. Behind one door is a car, behind both other doors a goat. You choose one of the doors. Monty Hall opens one of the other two doors, and shows that there is a goat behind it. You are now allowed to switch to the other door that is still closed. Is it smart to switch?

The Monty Hall Wikipedia Wars

(Gill 11, Mlodinow 08)

- Both sides **agree**:
 1. It is better to switch!
 2. To model problem correctly, you must take Monty's Protocol into account – what does Monty do when he has a choice?

The Monty Hall Wikipedia Wars

(Gill 11, Mlodinow 08)

- Both sides **agree**:
 1. It is better to switch!
 2. To model problem correctly, you must take Monty's Protocol into account – what does Monty do when he has a choice?
- “war” is about how to **prove** that switching is better:
 - “strictly Bayesian”: via conditioning, with additional assumption that Monty chooses by tossing a fair coin
 - credal set (imprecise probability, ambiguity)-based: make no assumptions on Monty and show e.g. that switching is minimax optimal

The Monty Hall Wikipedia Wars

(Gill 11, Mlodinow 08)

- Both sides **agree**:
 1. It is better to switch!
 2. To model problem correctly, you must take Monty's Protocol into account – what does Monty do when he has a choice?
- “war” is about how to **prove** that switching is better:
 - “strictly Bayesian”: via conditioning, with additional assumption that Monty chooses by tossing a fair coin
 - credal set (imprecise probability, ambiguity)-based: make no assumptions on Monty and show switching is e.g. dominating strategy

The Model on which they agree

- Suppose Contestant invariably chooses door a .
- Let RV Y denote location of car: $Y \in \{a, b, c\}$
- Let RV X denote Monty's action:

$$X \in \{ \text{open}(b), \text{open}(c) \}$$

$$X = \text{open}(c) \text{ means Monty opens door } c.$$

$$X = \text{open}(b) \text{ means Monty opens door } b.$$

This is really what Gilboa called the ‘eternal discussion’

The Model on which they agree

- Suppose Contestant invariably chooses door a .
- Let RV Y denote location of car.
- Let RV X denote Monty's action.

$$P(Y = a) = P(Y = b) = P(Y = c) = \frac{1}{3}$$

$$P(X = \text{open}(b) | Y = b) = P(X = \text{open}(c) | Y = c) = 0$$

The Point Probabilists' Side

- Suppose Contestant invariably chooses door a .
- Let RV Y denote location of car.
- Let RV X denote Monty's action.

$$P(Y = a) = P(Y = b) = P(Y = c) = \frac{1}{3}$$

$$P(X = \text{open}(b) | Y = b) = P(X = \text{open}(c) | Y = c) = 0$$

$$P(X = \text{open}(b) | Y = a) = P(X = \text{open}(c) | Y = a) = \frac{1}{2}$$

The Sets-of-Probabilities Side

- Suppose Contestant invariably chooses door a .
- Let RV Y denote location of car.
- Let RV X denote Monty's action.

$$P(Y = a) = P(Y = b) = P(Y = c) = \frac{1}{3}$$

$$P(X = \text{open}(b) | Y = b) = P(X = \text{open}(c) | Y = c) = 0$$

$$P(X = \text{open}(b) | Y = a) = 1 - P(X = \text{open}(c) | Y = a) \in [0, 1]$$



Dilation

$P(X = \text{open}(b) | Y = a) \in [0, 1]$ leads to

$$P(Y = a | X = \text{open}(b)) \in \left[0, \frac{1}{2}\right]$$

Instance of (partial) **dilation** (Seidenfeld, Wasserman 93):

Before observing X we had precise probability $P(Y = a)$
after we only know $P(Y = a | X = x)$ is in large superset

Dilation

$P(X = \text{open}(b) | Y = a) \in [0, 1]$ leads to

$$P(Y = a | X = \text{open}(b)) \in \left[0, \frac{1}{2}\right]$$

Instance of (partial) **dilation** (Seidenfeld, Wasserman 93):

Before observing X we had precise probability $P(Y = a)$
after we only know $P(Y = a | X = x)$ is in large superset

**“extra information \Rightarrow less knowledge
 no matter what you observe!”**

Assuming an Unbiased Monty

- To avoid dilation, tempting to become precise probabilist and **pretend that choices in protocol were made by fair coin tosses:**

$$P(X = \text{open}(b) | Y = a) = P(X = \text{open}(c) | Y = a) = \frac{1}{2}$$

...implying the familiar result

$$P(Y = a | X = \text{open}(b)) = \frac{1}{3}$$

$$P(Y = c | X = \text{open}(b)) = \frac{2}{3}$$

Assuming an Unbiased Monty...

- ..., i.e. use

$$P^o(Y = b \mid X = \text{open}(c)) := 2/3$$

is

1. "safe" and
2. minimax optimal

...under all **symmetric** decision problems and all, **even asymmetric Kelly gambling** problems

Safety for Decision Problems

Unbiased Monty is

1. "safe" under all symmetric loss functions: for all $P^* \in \mathcal{P}^*$:

$$E_{P^o}[\text{Loss}(Y, \delta_{P^o}(X))] = E_{P^*}[\text{Loss}(Y, \delta_{P^o}(X))]$$

where

\mathcal{A} = set of actions

Loss : $\mathcal{Y} \times \mathcal{A} \rightarrow \mathbb{R}$

$\delta_{P^o}(x) := \arg \min_{q \in \mathcal{A}} E_{P^o|X=x}[\text{loss}(Y, q)] = \text{Bayes act rel. to } P^o$

Safety

Unbiased Monty is

1. "safe" under all symmetric loss functions: for all $P^* \in \mathcal{P}^*$:

$$E_{P^o}[\text{Loss}(Y, \delta_{P^o}(X))] = E_{P^*}[\text{Loss}(Y, \delta_{P^o}(X))]$$

Example:

$\mathcal{A} = \{a, b, c\}$

Loss : $\mathcal{Y} \times \mathcal{A} \rightarrow \{0, 1\}$

Loss(Y, \hat{y}) = $1_{Y \neq \hat{y}}$

$\delta_{P^o}(\text{open}(c)) = b$; $\delta_{P^o}(\text{open}(b)) = c$.

$$= 1/3$$

Safety

Unbiased Monty is

1. "safe" under all symmetric loss functions: for all $P^* \in \mathcal{P}^*$:

$$E_{P^o}[\text{Loss}(Y, \delta_{P^o}(X))] = E_{P^*}[\text{Loss}(Y, \delta_{P^o}(X))]$$

Decision-Maker's pragmatic distribution

Bayes act based on P^o

'true' distribution

credal set

Safety

Unbiased Monty is

1. "safe" under all symmetric loss functions: for all $P^* \in \mathcal{P}^*$:

$$E_{P^o}[\text{Loss}(Y, \delta_{P^o}(X))] = E_{P^*}[\text{Loss}(Y, \delta_{P^o}(X))]$$

Second Example: **logarithmic scoring rule**

\mathcal{A} = set of prob. mass fn. on $\{a, b, c\}$

Loss : $\mathcal{Y} \times \mathcal{A} \rightarrow [0, \infty]$

Loss(Y, q) = $-\log q(Y)$

$\delta_{P^o}(\text{open}(c)) = (\frac{1}{3}, \frac{2}{3}, 0)$; $\delta_{P^o}(\text{open}(b)) = (\frac{1}{3}, 0, \frac{2}{3})$

$$= H(1/3)$$

What about nonsymmetric losses?

- 'asymmetric' means e.g. that if the car is behind door B, it is a **Ferrari**; if it is behind door C, it is a **Fiat Panda**
- Now pretending that Monty chooses by tossing a fair coin is **neither safe nor minimax optimal!**
- Except for asymmetric versions of log-loss! Then fair-coin assumption is still both **safe and minimax optimal!**

$$\text{Loss}(Y, q) = -\log \frac{q(Y)}{a(Y)}$$

Assuming an Unbiased Monty...

- ..., i.e. use

$$P^\circ(Y = b \mid X = \text{open}(c)) := 2/3$$

is

1. “safe”
2. minimax optimal
3. admissible

hence **PRETTY ADEQUATE**

....under all **symmetric** decision problems and all, **even asymmetric Kelly gambling** problems

Unbiased Monty

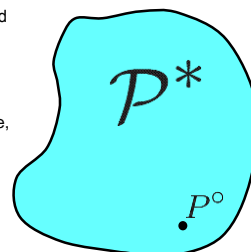
- Straightforward imprecise probability gives dilation
- Straightforward subjective Bayes is problematic for me: *why* would Monty be unbiased??
- Safe Probability Approach: *if* you are willing to make some assumption about loss function, it is *safe* to assume that Monty tosses a fair coin

Dependence on Task

- Safe Probability Approach: *if* you are willing to make some assumption about loss function, it is *safe* to assume that Monty tosses a fair coin
- This means that if you are told that the loss function is asymmetric, you may want to **change** your distribution
 - Similarly, if you're told in dilation problem that the probability that you have to make a prediction depends on V , you don't want to ignore V any more
 - Similarly, if, in 'objective Bayes', you change the sampling plan, you want to change the prior
 - **This is the price we pay for cooking a Bayesian omelet with imprecise eggs**

Conclusion: Towards A Theory of “Safe Probability”

- Compromise between 'strict' Bayes and imprecise probability theory
- \mathcal{P}^* has unique representative P° as in Minimum Description Length, **objective Bayes**, fiducial inference, 'MaxEnt...
- One absolutely crucial difference: we **restrict** use of P° to subset of all possible prediction tasks; eqv.



we 'condition' on the task