# STAT 530: Partial Identification

Mar. 22, 2010

---

## Deceptively simple problem

$X \sim$ Bernoulli($r$) in population.
Infer $r$ from random sample.

Slight twist: $X$ not measured well
(e.g. perhaps some subjects give wrong answer on questionnaire)

Measure $X^*$ instead of $X$ on sampled subjects, where
$SN = Pr(X^* = 1 | X = 1)$
$SP = Pr(X^* = 0 | X = 0)$

---

## Still looks relatively simple, presuming some info. on misclassification rates

$X^* \sim$ Bernoulli$\{rSN + (1 - r)(1 - SP)\}$ in population.
Say decide on prior:

$$p(r, SN, SP) \quad \propto \quad I_{(0,1)}(r) I_{(a,1)}(SN) I_{(b,1)}(SP)$$

### Two issues

- MCMC and parameterization
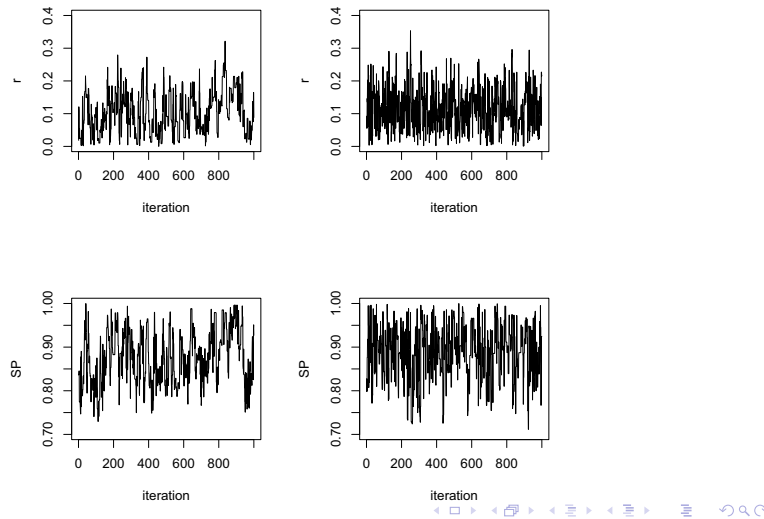- Information flow

---

## MCMC

### Example

- Data: 20/100 have $X^* = 1$
- Prior information: $a = 0.9$, $b = 0.7$

1. Univariate RWMH in original parameterization $(r, SN, SP)$, tuned to have approx. 50% acceptance rate for each component.

2. Univariate RWMH in new parameterization $(\tilde{r}, SN, SP)$, tuned to have approx. 50% acceptance rate for each component.
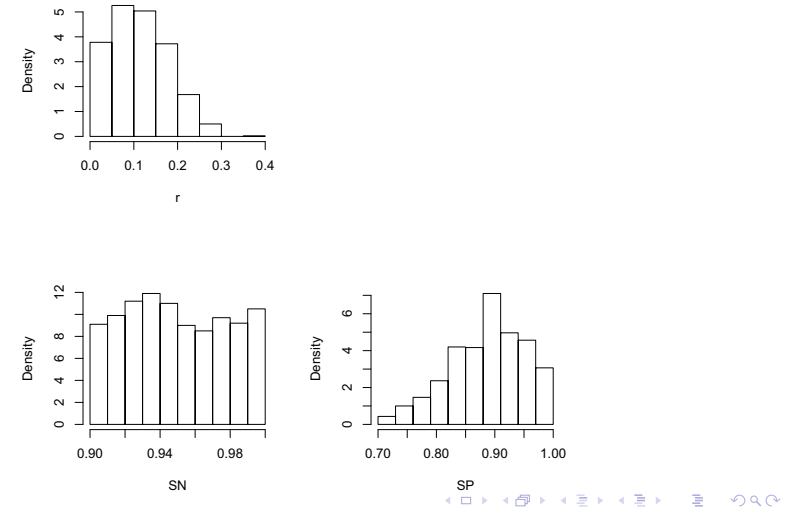
Intuition for why 2 might work better?

## Traceplots



## Information flow:
## Where does the 'extra-prior' info. about $SP$ come from???



## Intuition

Whereas $r$ and $(SN, SP)$ are independent *a priori*,
$\tilde{r}$ and $(SN, SP)$ are **dependent** *a priori*

Thus we learn about $\tilde{r}$ directly from the data,
and the prior dependence then implies something about $(SN, SP)$

[Also bear in mind that the quantity of most interest, $r$, can be
regarded as a function of $(\tilde{r}, SN, SP)$.]

## More formally

As $n \to \infty$

$p(\tilde{r}|\text{Data}) \to$

$p(SN, SP|\text{Data}) \to$

# We have just seen a simple example of a **partially identified** model

Arise somewhat generally when data are imperfect, but one isn't sure to what extent.

Characterized by the large-sample limit of the posterior distribution on the parameter of interest being:

- narrower than the prior
- wider than a single point

Point estimators (e.g. posterior mean) are necessarily biased.

But interval estimators (e.g. credible interval) reflect this appropriately.