Bayesian Inference, continued

In the example of polling from last time, there is in fact some prior knowledge - for example, polls in other states show Clinton and Sanders close to even. We have the notion that Illinois Democrats are not that different than Democrats in other states.

How to make use of this knowledge? With a richer class of priors:

We took $f_Y$ as Beta, so

$$f_Y(y) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} y^{\alpha-1}(1-y)^{\beta-1}$$

(remember: the Uniform dist is a special case of the Beta dist when $\alpha = \beta = 1$).

It turns out that as $\alpha$ and $\beta$ get bigger, and are not too different, the Beta distribution begins to resemble a Normal dist.

For $N(\mu, \sigma)$, as the

$$P(|Y - \mu| < \sigma) \approx \tfrac{2}{3} \approx .667$$

For Beta,

| $\alpha$ | $\beta$ | $P(|Y - \mu| < \sigma)$ |
|---|---|---|
| 1 | 1 | .577 |
| 2 | 2 | .626 |
| 3 | 3 | .644 |
| 5 | 5 | .659 |
| 4 | 6 | .661 |
| 1 | 19 | .812 ← SKEWED! |
| 10 | 10 | .671 |

...etc.

So suppose we expect that the vote is approximately split, with $Pr(.4 < Y < .6) \approx \tfrac{2}{3}$.

So
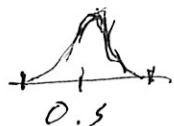
$$E(Y) \approx 0.5$$
$$Var(Y) \approx 0.1$$

②

For Beta distributions,
$$E(Y) = \frac{\alpha}{\alpha+\beta}, \quad Var(Y) = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$$
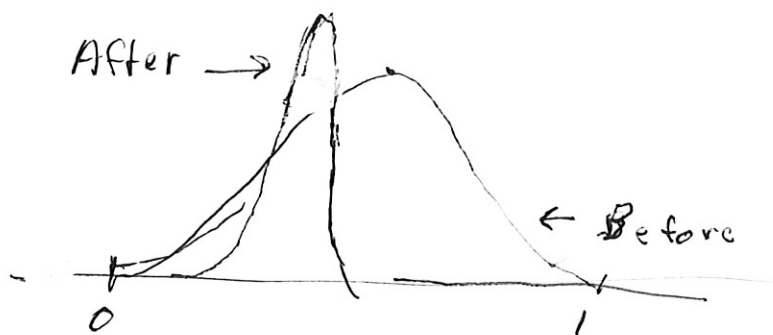
Set these equal to $0.5$ and $(0.1)^2$, solve to get $\alpha = \beta = 12$



$$f_Y(y) = \frac{\Gamma(24)}{\Gamma(12)\Gamma(12)} y^{11}(1-y)^{11}$$

has $E(Y) = 0.5$, $\sqrt{Var(Y)} = 0.1$

$$f(y/x) \propto f_Y(y) \, p(x/y)$$
$$= (\text{constant}) \, y^{11}(1-y)^{11}\binom{100}{40} y^{40}(1-y)^{60}$$
$$\propto y^{51}(1-y)^{71}$$



After →
← Before
0
1

[ Beta with $\alpha = 52$, $\beta = 72$ ]

Before: $E(Y) = 0.5$

After: $E(Y/X=40) = \frac{52}{124} = 0.42$

③

For Beta distributions,

$$E(Y) = \frac{\alpha}{\alpha+\beta} = \mu_Y$$

$$Var(Y) = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)} = \sigma_Y^2$$

$$\sigma_Y^2 = \overset{=\mu_Y}{\left(\frac{\alpha}{(\alpha+\beta)}\right)\left(\frac{\beta}{(\alpha+\beta)}\right)}\frac{1}{(\alpha+\beta+1)}$$

$$\frac{\beta}{\alpha+\beta} = \frac{\beta+\alpha-\alpha}{\alpha+\beta} = \frac{\beta+\alpha}{\alpha+\beta} - \frac{\alpha}{\alpha+\beta}$$

$$= 1 - \mu_Y$$

so

$$Var(Y) = \frac{\mu_Y(1-\mu_Y)}{\alpha+\beta-1}$$

$$\mu_Y = 0.5 \qquad \sigma_Y^2 = (0.1)^2$$

So:
$$(0.1)^2 = \frac{(0.5)(1-0.5)}{\alpha+\beta+1}$$

$$\alpha+\beta+1 = \frac{(0.5)^2}{(0.1)^2} = 25; \quad \alpha+\beta = 24$$

but

$$\mu_Y = \frac{\alpha}{\alpha+\beta} = 0.5 \implies \frac{\alpha}{24} = 0.5 \implies \alpha=12, \beta=12$$

(3A)

Note: We can interpret the posterior expectation as a weighted average:

$$\frac{\alpha + x}{\alpha + \beta + n} = \left(\frac{\alpha + \beta}{\alpha + \beta + n}\right) \cdot \frac{\alpha}{\alpha + \beta} + \left(\frac{n}{\alpha + \beta + n}\right) \cdot \frac{x}{n}$$

weights add to 1    prior expectation    sample fraction

Case 1:  $\alpha + \beta$  large relative to $n$
  ("strong prior information")

Then $\frac{\alpha + \beta}{\alpha + \beta + n} \approx 1$,  $\frac{n}{\alpha + \beta + n} \approx 0$

Case 2:  $n$ large relative to $\alpha + \beta$
  ("weak prior information")

Then $\frac{\alpha + \beta}{\alpha + \beta + n} \approx 0$,  $\frac{n}{\alpha + \beta + n} \approx 1$

Case 1:  $\frac{\alpha}{\alpha + \beta}$      Case 2:  $\frac{x}{n}$

otherwise,  a compromise!

$(4)$

# Bayes for Normal

$\theta$ is the **true value**.

We take the **prior** $f(\theta)$

$$f(\theta) = \frac{1}{\sqrt{2\pi}\,\sigma} e^{-\frac{1}{2}\left(\frac{\theta-\mu}{\sigma}\right)^2}, \quad -\infty < \theta < \infty$$

$$E(\theta) = \mu \qquad Var(\theta) = \sigma^2$$

$X$ is the **observed** **value** with error $\sim N(0, \gamma^2)$ So

$$X = \theta + error \text{ is } N(\theta, \gamma^2)$$

$$f(x|\theta) = \frac{1}{\sqrt{2\pi}\,\sigma} e^{-\frac{1}{2}\left(\frac{x-\theta}{\gamma}\right)^2}, \quad -\infty < x < \infty$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{likelihood}$$

$f(\theta|x)$, the **posterior** will be $N(A, B^2)$

[squares completed in Stigler, chapter 4]

$$A = \frac{\gamma^2 \mu + \sigma^2 x}{\gamma^2 + \sigma^2} \qquad\qquad B^2 = \frac{\gamma^2 \sigma^2}{\gamma^2 + \sigma^2}$$

$$\underbrace{\qquad\qquad}_{\substack{\text{Weighted Average} \\ \text{of } \mu \text{ and } X}} \qquad\qquad \underbrace{\qquad\qquad}_{\substack{\text{a posteriori} \\ \text{uncertainty}}}$$

(5)

## Example:

Measure a weight with an imperfect scale:

Scale makes errors with std. deviation 1 kg., normally distributed:

Distribution of scale readings $\longrightarrow$

standard normal, $\mu = y$

true weight

$X$ = recorded weight

$Y$ = true weight

$f(x|y) \sim \mathcal{N}(y, 1)$

$f_Y(y) \longrightarrow ??$

$f_Y(y) \searrow$

40    100    140

Say $\mathcal{N}(\mu, \sigma^2)$

$\mu = 100$ kg

$\sigma^2 = (10)^2 = 100$

[Why? Maybe have rough idea, say from number of people needed to lift it]

( A priori: $P(90 < Y \leq 110)$ 
$P(|Y - 100| \leq 10)$ $\Big\}$ $\approx \dfrac{2}{3}$

⑥

Given "data" $x$, want $f(y/x)$.

$$f(y/x) \propto f_Y(y) f(x/y)$$

$$= \frac{1}{\sqrt{2\pi}\,\sigma} e^{-\frac{1}{2}\left(\frac{y-\mu}{\sigma}\right)^2} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-y)^2}$$

$$\propto e^{-\frac{1}{2}\frac{(y-\mu)^2}{\sigma^2} - \frac{1}{2}(x-y)^2}$$

messy completion of squares not shown

$$\begin{cases} = e^{-\frac{1}{2\sigma^2}\left[(y-\mu)^2 + \sigma^2(x-y)^2\right]} \\ \propto e^{-\frac{1}{2}\frac{(Y-A)^2}{B}} \qquad \longleftarrow \end{cases}$$

functional form, aka "Business Part"

$$A = \frac{x\sigma^2}{\sigma^2+1} + \frac{\mu \cdot 1}{\sigma^2+1}$$

weighted average of $x$, $\mu$ and $x$.

$$B = \frac{\sigma^2}{(\sigma^2+1)}$$

$f(y/x)$ is $N(A, B)$

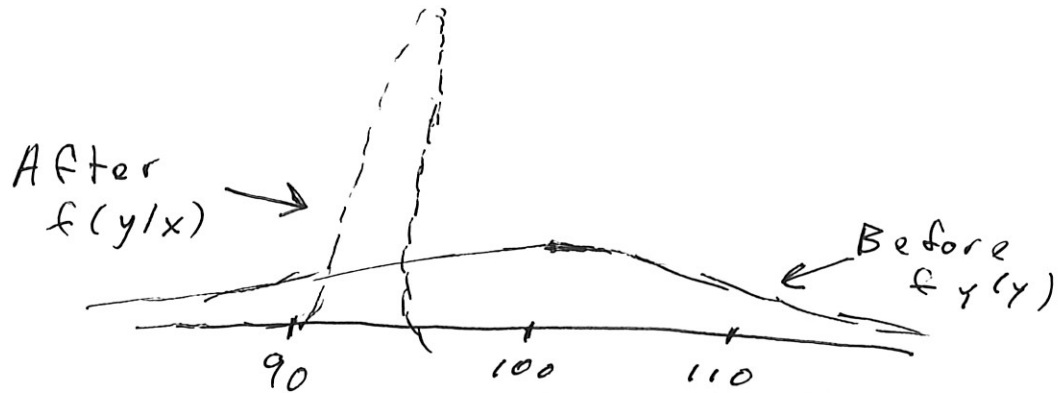$E(Y/x) = A$    (between $x$ and $\mu$)

If $\sigma^2$ small (much prior info)

    $A$ near $\mu$

If $\sigma^2$ large (little prior info)

    $A$ near $x$

After
f(y/x) →

Before
← $f_Y(y)$



90    100    110

Case for $\sigma^2 = (10)^2$, $\mu = 100$

$X = 90$

That is:

$$f_Y(y) \qquad N(100, 10^2)$$

$$f(x/y) \qquad N(y, 1)$$

$$A = \frac{100}{101} \cdot X + \frac{1}{101} \cdot \mu = 90.9$$

$$B = \frac{100}{101}$$

$$f(y/x) \qquad N\left(90.9, \frac{100}{101}\right)$$

Bayes's Theorem Processes
Information!

⑧

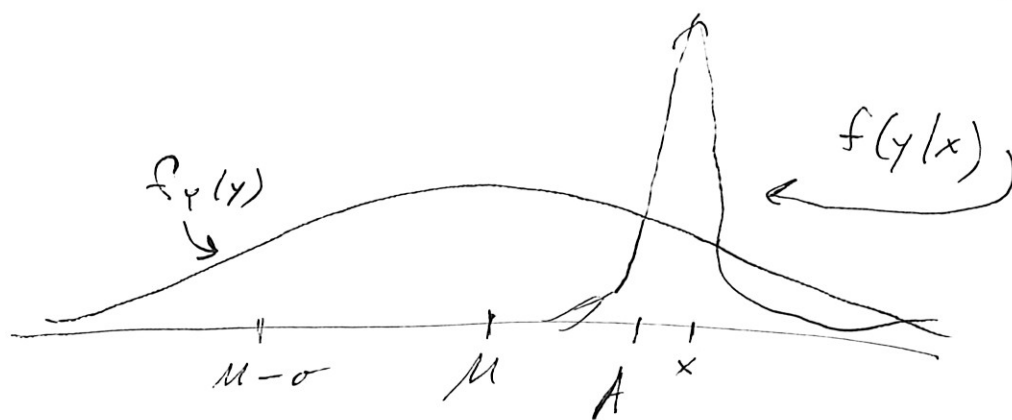In summary, for the normal dist, we ha

$$f(y) \qquad N(\mu, \sigma^2)$$
$$f(x|y) \qquad N(y, 1)$$
$$f(y|x) \qquad N(A, B)$$

$$A = x \cdot \frac{\sigma^2}{\sigma^2 + 1} + \mu \cdot \frac{1}{\sigma^2 + 1}$$

$$B = \frac{\sigma^2}{\sigma^2 + 1}$$

$$\left( \text{so if } \lambda = \frac{\sigma^2}{\sigma^2 + 1} , \quad A = x \cdot \lambda + \mu(1 - \lambda) \right)$$



$f_Y(y)$

$f(y|x)$

$\mu - \sigma \qquad \mu \qquad A \qquad x$

# Intro to Maximum Likelihood

Today we'll go further in considering Statistical Inference.

Recall that we would like to know the "state of nature" $\theta$. More exactly $\theta$ is a parameter that represents such a state.

We will learn about $\theta$ by considering $X$ $[= (x_1, ..., x_n)]$, the data.

We need a model to describe the relation between $\theta$ and $X$.

Specifically, $X$, given $\theta$, is a random variable with distribution $p(x|\theta)$ or $f(x|\theta)$.

Ex: $\theta = $ fraction of votes
$X = $ # of 100 sampled
$$p(x|\theta) = \binom{100}{x} \theta^x (1-\theta)^{100-x}$$
$x = 1, 2, ..., 100$

Ex: $\theta = $ true weight
$X = $ what scale says
$f(x|\theta)$ $N(\theta, 1)$

Ideal Goal: Find $f(\theta/x)$.
   ie: After we have data
("given data"), we want to know
the probability of various values
of $\theta$.

So far, we've used
        Bayes's Theorem:
        $f(\theta/x) \propto \underbrace{f(\theta)}_{prior} f(x/\theta)$

   Gives what we want. But:
      it requires $f(\theta)$.
      $f(\theta)$ is controversial —
                  How to get it?
                  what does it mean?
                  subjective bias — disagreements

   OK. How about a more limited goal?
      We won't use $f(\theta)$.
      Instead, we will work only
         with $f(x/\theta)$. Then we'll

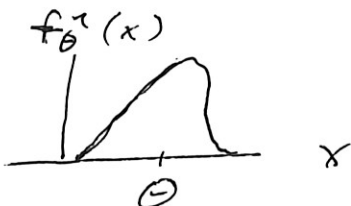Estimate a Point, not a Distribution
_____
      We treat $\theta$ as fixed (a 'given')  X
as  random. We want an estimate of
$\hat{\theta} = f(x)$ that is likely to be close
to $\theta$.

        ⑪

$\hat{\theta}$ depends on $X$

$\hat{\theta}$ is short for $\hat{\theta}(x)$
(or $\hat{\theta}(x_1, \ldots x_n)$)

$X$ is a random variable, so
$\hat{\theta}$ is a random variable
_____

What does "$\hat{\theta}$ likely to be near $\theta$" mean? From our "given $\theta$" perspective, $\hat{\theta}$ has a distribution $f_{\hat{\theta}}(x)$ or $f_{\hat{\theta}}(x|\theta)$:



We want this distribution to be concentrated near and/or centered at $\theta$.

Defs: $\hat{\theta}$ is <u>unbiased</u> if $E(\hat{\theta}) = \theta$, whatever $\theta$ is (i.e. $\int_{-\infty}^{\infty} x \, f_{\hat{\theta}}(x|\theta) \, dx = \theta$ for all $\theta$).

<u>Bias</u> $= E(\hat{\theta}) - \theta$

<u>Mean Error</u> $= E(|\hat{\theta} - \theta|)$

<u>Mean Square Error</u> $= E[(\hat{\theta} - \theta)^2]$
("MSE")

(13)

It turns out the MSE has a particularly clear interpretation:

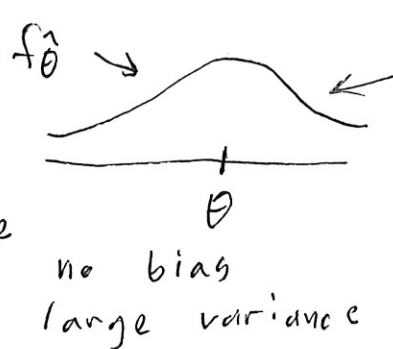$$MSE_{\hat{\theta}}(\hat{\theta}) = E(\hat{\theta} - \theta)^2$$

$$= E[(\hat{\theta} - E(\hat{\theta})) + (E(\hat{\theta}) - \theta)]^2$$

$$= E[(\hat{\theta} - E(\hat{\theta}))^2 + 2(\hat{\theta} - E(\theta))(E(\theta) - \theta) + (E(\hat{\theta}) - \theta)^2]$$

$$= E[\hat{\theta} - E(\hat{\theta})]^2 + 2(E(\hat{\theta}) - \theta)\underline{E(\hat{\theta} - E(\hat{\theta}))}$$

$$+ (E(\hat{\theta}) - \theta)^2$$

$Var(\hat{\theta}|\theta)$

$(B(\theta))^2$
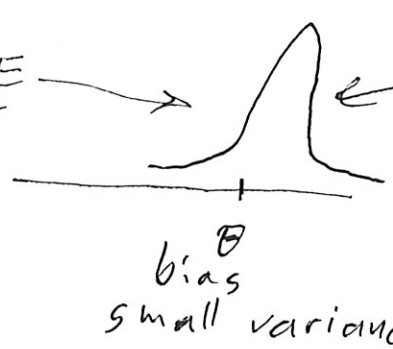
$= E(\hat{\theta}) - E(\hat{\theta}) = 0$

Hence

$$MSE(\hat{\theta}) = Var(\hat{\theta}) + (Bias)^2$$

$\uparrow$      $\uparrow$      $\uparrow$

"expected error"    error from spread of data points    error from bias

Tradeoff between bias and variance

Case I   $f_{\hat{\theta}}$

$\leftarrow$ SAME MSE $\rightarrow$

Case II   $f_{\tilde{\theta}}$

$\theta$

no bias large variance

$\theta$

bias small variance



(13)

Example : X Binomial $(n, \theta)$

$$\hat{\theta} = \frac{X}{n} . \quad E(\hat{\theta}) = \frac{E(x)}{n} = \frac{n\theta}{n} = \theta$$

$$\underline{unbiased}$$

Example : Same dist, but

$$\hat{\theta}^* = \frac{X+1}{n+2}$$

This estimator for $\hat{\theta}$ is what we'd get in a Bayesian analysis with $f(\theta)$ uniform on $[0,1]$. Then the posterior dist. would be Beta $(x+1, n-x+1)$ with $E(\theta/X=x) = \frac{X+1}{n+2}$. We are not being Bayesian here, but we can still use the estimator.

$$E(\hat{\theta}^*) = \frac{E(X)+1}{n+2} = \frac{n\theta + 1}{n+2} \neq \theta !$$

$$\underline{Biased} !$$