

Q: How do we know when a statistical model is 'good'?

Useful A: A

DD AM
discussion
section

model is 'good' if it makes accurate & well-calibrated predictions of data $\textcircled{1}$

values not used in fitting it with out-of-sample predictive validation

out-of-sample models, fitting means estimating its parameters & obtaining well-calibrated uncertainty assessments about (a) those parameters & (b) predictions of new data

To instantiate $\textcircled{+}$ we need a measure of goodness of fit of a model

single example

$$\left\{ \begin{array}{l} (\mu, \sigma) \sim \text{LI}(\mu, \sigma) \\ \{I_i | \mu, \sigma\} \stackrel{\text{i.i.d.}}{\sim} N(\mu, \sigma^2) \end{array} \right\} \textcircled{+}$$

$(i=1, \dots, n)$

for each moderate n we showed (when looking at Bayesian analyses of Gaussian sampling models) that the predictive distribution

for a new data value Y_{new} is $\boxed{\begin{matrix} Y = (y_1, \dots, y_n) \\ \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \end{matrix}}$
 $(Y_{\text{new}} | Y \text{ (SINK)} [PM: LI] \text{ (B)}) \sim N(\bar{y}, s_n^2)$

$$s_n = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2}$$

with a few batch of data

a natural goodness-of-fit measure is $Y_{\text{new}} = (y_{n+1}, \dots, y_{n+l})$

$$\frac{1}{l} \sum_{i=1}^l (y_{n+i} - \bar{y})^2 \triangleq \text{GoF}$$

this is the squared Euclidean distance between

the vector Y_{new} and the scalar \bar{y} , scaled by sample size to create a mean discrepancy but we can't compute this until Y_{new}

arises $\left\{ \begin{matrix} 2 \\ \text{solutions} \end{matrix} \right\}$ ① cross-validation | ② penalized complexity

Cross validation (CV) ①

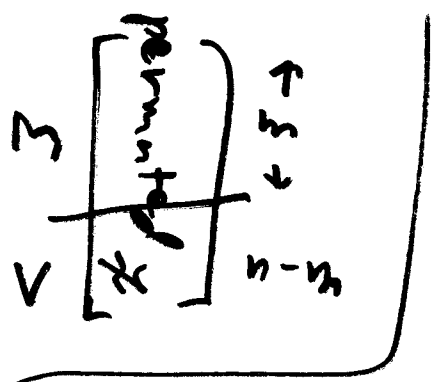
"1-way CV": use all of (cheating) y to assess predictive

③ $\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = y$
observed data vector

quality by calculating $\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y}_n)^2$

this is easy but may well understate our actual uncertainty about y_{new} , because we used the data twice: once to get \bar{y}_n and once to see how good \bar{y}_n is (this can be remedied with a complexity penalty [see ② below])

"2-way CV" (A) randomly partition



y into modeling (M) and validation (V) subsets of size m & $(n-m)$ (respectively)

for a choice of $1 \leq m \leq n-1$; (B) use the mean of the (M) subset values to predict the (V) subset values & compute the GOF summary on (V)

(C) repeat ^{(A), (B)} many times & average the results & summarize (4)

Q: How does this generalize to other models?
A: The "2 way CV" structure is already completely general for settings in which the dataset \mathcal{D} is a vector γ of real numbers.

But the natural goodness of fit measure in other models is less clear.

It turns out that looking at the log likelihood function is illuminating.

In model ~~(*)~~ (p. 1 bottom)

$$\text{This is } \ell(\mu, \sigma | \gamma [SM:N] B) = -\frac{n}{2} \log \sigma^2 - \frac{n}{2} \log 2\pi - \frac{1}{2\sigma^2} \sum_{i=1}^n (\gamma_i - \mu)^2$$

recall that the MLEs

in this model are $\hat{\mu}_{MLE} = \bar{\gamma} = \frac{1}{n} \sum_{i=1}^n \gamma_i$ and $\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (\gamma_i - \bar{\gamma})^2$, so the maximum log lik.

value is $\mathcal{L}(\hat{\mu}_{MLE} \hat{\sigma}_{MLE} | \mathcal{D} [SM: N] \mathcal{B}) =$

$$-\frac{n}{2} \log \left[\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \right] - \frac{n}{2} \log 2\pi - \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\frac{2}{n} \sum_{i=1}^n (y_i - \bar{y})^2}$$

$$= -\frac{n}{2} \left\{ \log \left[\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \right] - \log 2\pi - \frac{1}{2} \right\}$$

but this is the natural GOF measure for the Gaussian model (!)

This suggests

a general rule: use the maximum log likelihood value to define goodness of fit

in practice people use

$$-2 \cdot \mathcal{L}(\hat{\theta} | \mathcal{D} [SM] \mathcal{B})$$

as the

general GOF measure, with

small values
↓
good fit

⑥
 $-2 \ln(\hat{\theta} | D [SM] B)$ is called the
frequentist deviance of the model based on IID
realizations from the sampling model [SM]

The multiplication by (-2) is because of
a frequentist story under which it
turns out that $(-2 \ln_{\max})$ can be
calibrated by comparing it to values
of the χ^2 distribution (a special
case of the $\Gamma(a, b)$ family)

However,
the deviance is based on the "I may CV"
idea of using the entire dataset to
compute \ln_{\max} ; to avoid over-fitting

it turns out that we need to penalize the model leading to ℓ_{max} according to its complexity (otherwise we could get a superb value of ℓ_{max} just by ~~fitting the data perfectly~~ fitting the data perfectly)



model (polynomial of order $n=7$, $n = \# \text{ data points}$) with a superb value of ℓ_{max} but validates poorly out of sample (high complexity)

model (linear) that validates well out of sample, but has a small ℓ_{max} but low complexity

Important note

Q: Since the likelihood function is only defined up to a positive constant multiple, can't we just make ℓ_{max}

as large as we like just by choosing ⑧
not constant?

A: Yes; that's temporarily

embarrassing. The fix is as follows:

when computing the likelihood function for model comparison & assessment of GOF, (a) write down the joint sampling distribution with its correct normalizing constant; (b) define the likelihood function for GOF evaluation as equal to the function in (a), including the normalizing constant from (a) in this definition (I did this in the example above but didn't mention it until now).

3 math
impulses
following
development
of a good
idea

- ① try to break it, to find out its scope of validity;
- ② try to generalize it;
- ③ use it to unify ideas that were previously

regarded as separate

key statistical

unifying idea for parametric sampling models: the Exponential Family (EF)

(not to be confused with the Exponential sampling model, which is (confusingly) itself a member of

the EF)

In 1935-1936, 3 sets of researchers

(Cambridge)	(École Normale Supérieure)	(Columbia)
(E. J. G. Pitman & John Wishart)	(G. Darrois)	(B. O. Koopman)

independent by & more or less simultaneously published the same excellent idea, on which the EF is based (Lecture Notes part 6 pp. 149 →)

General Parametric Model for a data vector $\underline{y} = (y_1, \dots, y_n)$

sampling model

$$p(y_i | \underline{\theta} [SM] \mathcal{B})$$

marginal PMF/PDF for a single observation y_i .

$\underline{\theta} = (\theta_1, \dots, \theta_k)$

joint sampling model

$$(\underline{y}_i | \underline{\theta} [SM] \mathcal{B}) \stackrel{IID}{\sim}$$

$$p(\underline{y} | \underline{\theta} [SM] \mathcal{B}) =$$

$$p(y_i | \underline{\theta} [SM] \mathcal{B}) \quad (i=1, \dots, n)$$

$$\prod_{i=1}^n p(y_i | \underline{\theta} [SM] \mathcal{B})$$

this is a member of the EF

iff it can be written

$$f(\underline{y}) g(\underline{\theta}) \cdot \exp \left[\sum_{j=1}^k \phi_j(\underline{\theta}) \sum_{i=1}^n h_j(y_i) \right]$$

2 immediate advantages

① a set of sufficient statistics is then obvious: (minimal)

$$\left\{ \sum_{i=1}^n h_1(y_i), \dots, \sum_{i=1}^n h_k(y_i) \right\}$$

② This sampling model must have a conjugate prior, which is as follows:

$$p(\underline{\theta} | \text{[CP]} \mathcal{B}) = c g(\underline{\theta})^{\tau_0} \exp \left[\sum_{j=1}^k \phi_j(\underline{\theta}) \tau_j \right]$$

for some vector $\underline{\tau} = (\tau_0, \tau_1, \dots, \tau_k)$

example | $(I_i | \theta \mathcal{B}) \stackrel{\text{IID}}{\sim} \text{Bernoulli}(\theta)$

$k=1$ | $(0 < \theta < 1)$ $(i=1, \dots, n)$

$$\underline{\theta} = (\theta) \quad p(y_i | \theta \mathcal{B}) = \theta^{\tau_i} (1-\theta)^{1-\tau_i}$$

$$p(\mathbf{y} | \theta, \mathcal{B}) = \prod_{i=1}^n p(y_i | \theta, \mathcal{B}) = \theta^s (1-\theta)^{n-s} \quad (12)$$

$$\text{for } s = \sum_{i=1}^n y_i$$

to be in the EF club
this has to be expressible

$$\text{as } f(\mathbf{y})g(\theta) \exp\left[\phi_1(\theta) \sum_{i=1}^n h_1(y_i)\right]$$

($k=1$)

They really don't look the same,

but here's a good EF trick: for

$$\text{all } x > 0, \quad x = \exp[\log(x)]$$

✓

$$\theta^s (1-\theta)^{n-s} = \exp\left\{\log\left[\theta^s (1-\theta)^{n-s}\right]\right\}$$

$$= \exp\left[s \log \theta + \overset{(n-s)}{\cancel{n-s}} \log(1-\theta)\right]$$

$$= \exp\left[s \log \theta + n \log(1-\theta) - s \log(1-\theta)\right]$$

$$\theta^s (1-\theta)^{n-s} = \exp \left[n \log(1-\theta) + s \log\left(\frac{\theta}{1-\theta}\right) \right] \quad (13)$$

$$= \exp \left[n \log(1-\theta) \right] \cdot \exp \left[s \log\left(\frac{\theta}{1-\theta}\right) \right]$$

$$= \underbrace{1}_{f(\gamma)} \cdot \underbrace{\frac{(1-\theta)^n}{\theta^s}}_{g(\theta)} \exp \left[\underbrace{s}_{\sum_{i=1}^n h_1(\gamma_i)} \cdot \underbrace{\log\left(\frac{\theta}{1-\theta}\right)}_{\phi_1(\theta)} \right]$$

logistic regression

with $h_1(\gamma_i) = \gamma_i$ EF ✓

$\phi_1(\theta) = \log\left(\frac{\theta}{1-\theta}\right)$ is called the natural parameterization of the

Bernoulli(θ) sampling distribution

minimal
 the sufficient statistic is $s = \sum_{i=1}^n h_1(\gamma_i)$

and the conjugate prior, which we know is 14

$$p(\theta | [CP] B) = c \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

can be written as $c g(\tau) \tau_0 \exp[\phi(\theta) \tau_1]$:

$$= c (1-\theta)^{n \tau_0} \exp\left[\tau_1 \log\left(\frac{\theta}{1-\theta}\right)\right]$$

Some EF sampling models are regular

& some are not

Def. An EF

sampling model is regular if the

with
parameter
vector θ

support of the marginal
sampling dist of Σ_i does

not depend on θ

ex. ① $(Z_i | \mu, \sigma [SM: N] \mathcal{B}) \stackrel{IID}{\sim} N(\mu, \sigma^2)$ ⑩
 $(i=1, \dots, n)$

$\theta = (\mu, \sigma)$
 \sim
 $\sigma = \underline{\sigma}$
 \sim
 $\sigma^2(\mu, \sigma^2)$
 \sim
 $\sigma^2(\mu, \frac{1}{\sigma^2})$

$\text{Support}(Z_i) = \mathbb{R}$

does not depend on θ

so this SM is a regular member of the EF

ex. ②

$V(\hat{\mu}_{MLE}) = \underline{\underline{O}}(\frac{1}{n})$

$(Z_i | \theta [SM: n] \mathcal{B}) \stackrel{IID}{\sim} \text{Uniform}(0, \theta)$
 $(i=1, \dots, n)$ $(\theta > 0)$

$\text{Support}(Z_i) = (0, \theta)$ this depends on θ

this is a non-regular member of EF

$V(\hat{\theta}_{MLE}) = \underline{\underline{O}}(\frac{1}{n^2}) (!)$

Edwards, Lindeman, & Savage ~ 1940s (16)

psychometrics

LI prior:

low info =

flat =

diffuse =
(large spread)

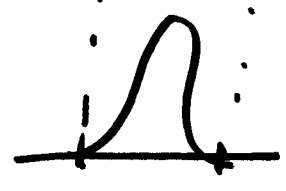
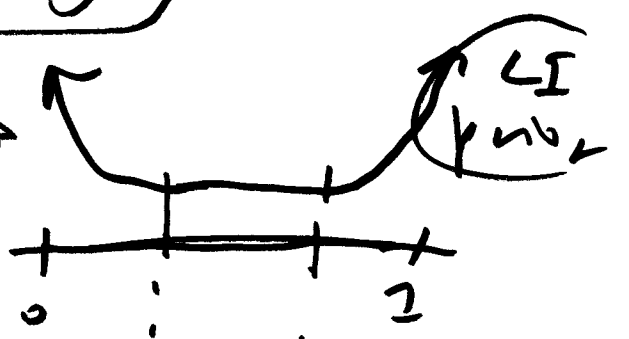
~~noninformative~~

$U(\text{left}, \text{right})$

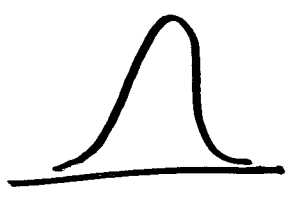
to obtain a LI prior, always choose 'left' & 'right' to avoid

~~at~~ Beta($\frac{1}{2}, \frac{1}{2}$)

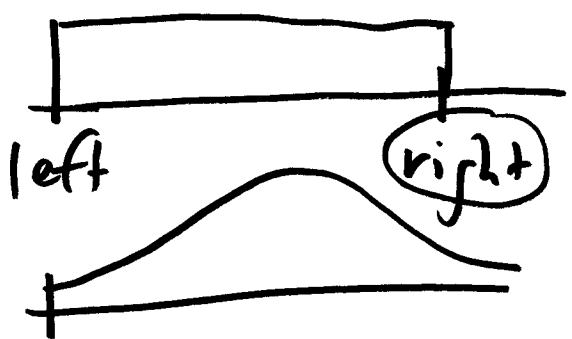
principle of stable estimation



likelihood

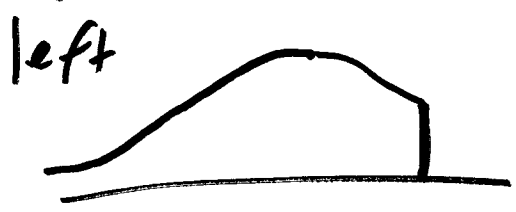


posterior =
lik



prior

likelihood



posterior

inappropriate truncation of lik. f's

$$p(\gamma_{\text{new}} | \gamma [SM] [PM] B) \quad \text{(hard)}$$

$k \sim_i \text{ in } [SM]$ (17)

$$= \int p(\gamma_{\text{new}}, \theta | \gamma [SM] [PM] B) p(\gamma_{\text{new}} | \theta [SM] [PM] B) d\theta$$

Sampling distribution for γ_{new}

(easy)

$$= \int p(\gamma_{\text{new}} | \theta [SM] [PM] B) d\theta$$

post. dist. $\int p(\theta | \gamma [SM] [PM] B) d\theta$

how to MCMC sample from this predictive distribution

$\textcircled{+}$ is actually

telling a mixture story in disguise:

$\textcircled{+}$ is a weighted average of [sampling distribution for given θ]

weighted by the post. dist. for θ

how to make random draw from Θ

① make a random draw of θ from the posterior, obtaining θ^*

② make a random draw from the sampling dist. using the θ^* value from ①

