

## validation of SBI (world)

- given test data  $\{X_i^{\wedge} \sim p(x)\}_{i=1}^N$  or  $\{X_i^{\wedge} = p(x | Y = \text{fixed})\}_{i=1}^N$   
 $\{X_i^* \sim p^*(x)\}_{i=1}^{N'}$  or  $\{X_i^* = p^*(x | Y = \text{fixed})\}_{i=1}^{N'}$

case (1)  $N' \approx N \gg 1 \Rightarrow$  compare the distribution of the samples  $\{X_i^{\wedge}, X_i^*\}$  using MMD or FID or density/coverage

case (2)  $N' = 0$  (no ground truth)  $\Rightarrow$  compare diversity of different approximations via Vendi score  
[if  $N' \gg 1$ , also compare against GT Vendi score]

case (3)  $N' = 1$  important case in practice

- in SBI, create GT via forward simulation,  $\{Y_i \sim p^s(Y), X_i = \phi(Y_i, \epsilon)\}$
- $\Rightarrow Y_1$  is a GT example for  $p(Y | X = X_1^{\wedge})$  with  $N' = 1$   
 $\uparrow$  fixed

• weather forecasts:  $p(Y = \text{rain tomorrow} | X = \text{weather up to now})$

"80% rain probability"

$Y_i = p^*(Y = \text{rain} | X = \text{weather so far}) \hat{=} \text{actual weather}$   
is a GT sample with  $N=1$

• idea: "calibration"

merge instances with same predicted confidence in a joint test set among all days with "80% rain prob." it should have rained in 80% of the cases

applied to classification.

$p(Y=k|x)$  is 80%, and answer is right

80% of the time  $\Rightarrow$  "well calibrated"

$> 80\%$  ———  $\Rightarrow$  "underconfident"

$< 80\%$  ———  $\Rightarrow$  "overconfident"

• realization: sort predicted sample  
(consider the inverse problem

$p(Y|x)$  with  $Y \in \mathbb{R}$

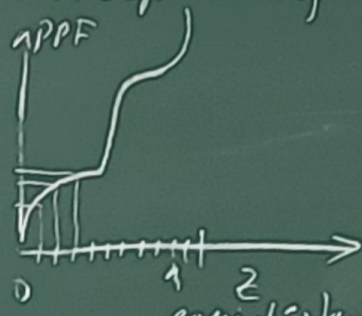
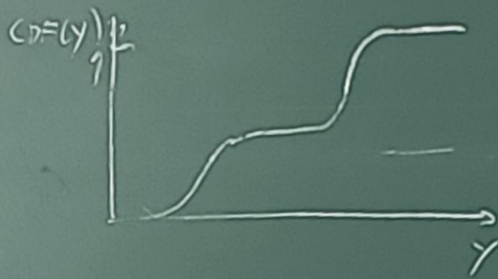
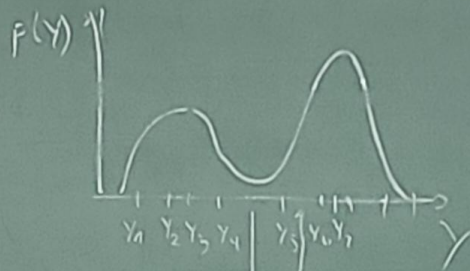
$\{Y_i = \hat{y}_i p(Y|x = \text{fixed})\}_{i=1}^N \cup \{Y_0 = Y^* \sim p^*(Y|x)\}$

sort into  $(Y_{[0]}, \dots, Y_{[N]})$

$\leftarrow$  sorted order

$Y_0 \rightsquigarrow Y_{[k]}$  calibration  $\hat{=} k \sim \text{uniform}(0, N)$

let  $p(Y)$  some 1-D prob,  $CDF(Y)$  corresponding cumulative dist. f.



equidistant quantiles

$$\frac{1}{(N+1)}, \frac{N}{(N+1)}$$

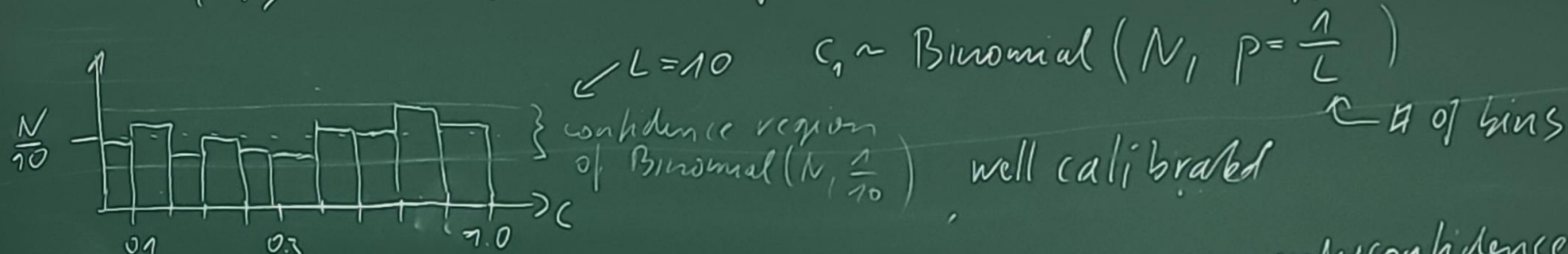
all intervals have same probability mass  
 $\Rightarrow$  the interval where a new sample lands  
 is uniformly distributed

- if quantiles are unknown, a sorted sample  $Y_{(i)} \sim p(Y)$  is a good approx.  
 probability mass in every interval is approximately equal

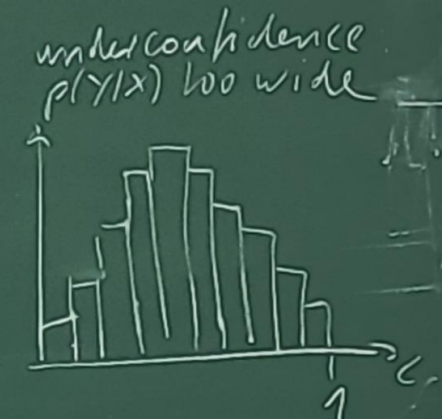
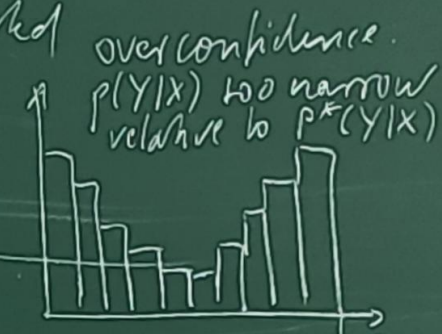
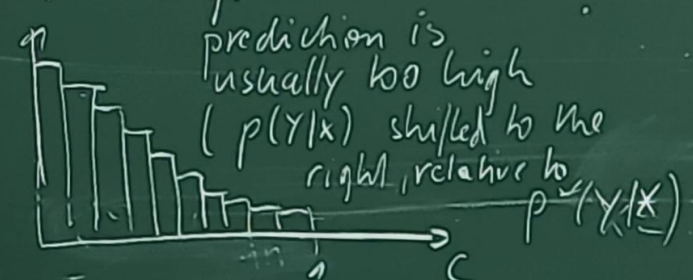
- alg. ① given GT  $Y_i^*$ , sample from model  $\{\hat{Y}_k \sim p(Y | X_i)\}_{k=1}^M$
- ② sort the joined set  $\{Y_{i0} = Y_i^*, \hat{Y}_{i1}, \dots, \hat{Y}_{iM}\} \Rightarrow (Y_{i[0]}, \dots, Y_{i[M]})$
- ③  $C_i = \frac{L_i}{M} \leftarrow$  index of  $Y_i^*$  in sorted order  $(C \in [0, 1])$

repeat this for many GT instances  $(X_i, Y_i) \Rightarrow$  sample  $\{C_i\}_{i=1}^N$

• evaluate  $\{C_t\}$  via histogram test: if model is calibrated (null hypothesis)



typical histograms when uncalibrated



• caveat: "well-calibrated" does not imply "accurate"  
a bad model that knows that it is bad, is still calibrated

toy example:  $t$  time, sample mean on day  $t$   $\mu_t \sim N(0, 1)$ , outcome  $Y_t \sim N(\mu_t, 1)$   
 $\rightarrow$  marginal distribution  $Y_t \sim p^*(Y_t) \sim N(0, \sigma^2 = 2)$

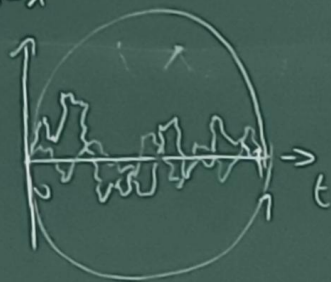
if model predicts  $Y_t \sim p(Y) \approx p^*(Y_t)$ , it is calibrated, but predicting  $X_t \sim p(X|\mu_t)$  is better  
 both are calibrated but variances differ

• alternative visualization via empirical CDF of  $\{C_i\}_{i=1}^N$   $C_i \in [0, 1]$

$$\text{CDF}(t) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}[C_i \leq t] \quad \text{define } v(t) = \text{CDF}(t) - t$$

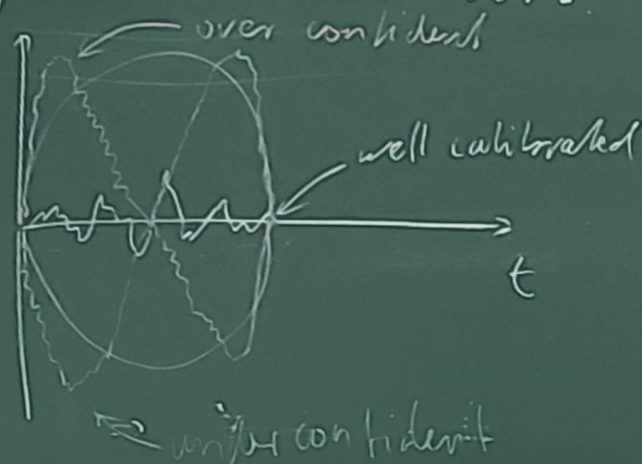
$\mathbb{Z} \text{CDF}^*(t)$

draw  $v(t)$



- one curve per element of  $Y \in \mathbb{R}^D$
- if null hypothesis "well calibrated" is true curves are "Brownian bridges"  $\equiv$  random walk with fixed start and end point

typical calibration errors



- confidence intervals of Brownian bridge are analytically known (in practice, draw

$$\text{Stdev}(t) = \sqrt{\frac{t(1-t)}{N}}$$

$$\frac{1}{\sqrt{N}} v(t) \quad \text{with } \text{stdev}(t) = \sqrt{t(1-t)}$$

$\mathbb{Z} \leftarrow$  implies 95% confidence region

joint calibration checks: instead of testing  $Y_i = (Y_{i1}, \dots, Y_{iD})$  on feature at a time, check the entire vector

$\Rightarrow$  reduce problem to 1-D via "energy" or "surprisal" distribution

$$e(Y | X_i) = -\log p(Y | X_i) \quad \left. \begin{array}{l} \text{energy if } p(Y) = \frac{1}{Z} \exp(-e(Y)/kT) \\ \text{(Gibbs distribution)} \end{array} \right\}$$

define energy distribution via  
"how often do we see  $e(Y) = E$ "

$$-\log p(Y) = \frac{e(Y)}{kT} + \log Z$$

$$p(E | X_i) = \int_Y \delta(e(Y | X_i) - E) p(Y | X_i) dY$$

example if  $Y \sim N(0, \Sigma_D)$ , then  $-\log p(Y) = \frac{\|Y\|^2}{2} + \text{const}$

$$p(E) = p\left(\frac{\|Y\|^2}{2} = E\right) = 2 \chi_D^2(2E)$$

$$E = \text{const} = \left\{ Y \mid \frac{\|Y\|^2}{2} = E, \text{ i.e. surface} \right.$$

of a sphere with radius  $\sqrt{2E}$

alg: (1) for  $i=1, \dots, N$

(a)  $\{ \hat{Y}_{ik} \sim p(Y | X_i) \}_{k=1}^M$  and  $Y_i^*$  the GT from test set

(b)  $e_{ik} = -\log(\hat{Y}_{ik} | X_i)$ ,  $e_{i0} = -\log(Y_i^* | X_i)$

(c) sort into  $(e_{i(0)}, \dots, e_{i(M)})$ , let  $[k]$  index of  $Y_i^*$  after sorting

(d) define  $c_i = \frac{[k]}{M}$

(2) use histogram or CDF methods to analyse distribution of  $\{c_i\}_{i=1}^N$

---

weather forecast  $p(Y = \text{rain} | X = \text{weather solar, date})$  Tilman Gneiting

model ignores weather solar  $\Rightarrow$  makes predictions only according to date  
prediction is calibrated, if  $p(Y | \text{date})$  is the correct statistics for "date"  
but much less accurate if it exploited weather correlation  
between subsequent days