# Agenda

- **Day-#6**
  - Tutorial part 1 (Bayes Rule with Gaussians)
    - 3 Exercises + 1 Bonus
  - Tutorial part 2 (Causal inference + Gaussians)
    - 1 Exercise
  - Tutorial part 3 (Fitting to data)
    - 7 Exercises
  - Tutorial part 4 - Bonus (Bayesian decision theory)
    - 2 Exercises

**W2D1_pod 031**

# Objective

Developing a <u>Bayesian model</u> for localizing sounds based on audio and visual cues.

This model will combine <u>**prior** information</u> about where sounds generally originate with sensory information about the <u>**likelihood**</u> that a specific sound came from a particular location.

Resulting <u>**posterior distribution**</u> not only allows us to make optimal decision about the sound's origin, but also lets us quantify how uncertain that decision is.

Bayesian techniques are therefore useful <u>**normative models:**</u> the behavior of human or animal subjects can be compared against these models to determine how efficiently they make use of information.

Note: fundamental building blocks for Bayesian statistics: the Gaussian distribution and the Bayes Theorem.

**W2D1_pod 031**

# Central limit theorem

The **central limit theorem** states that if you have a population with mean $\mu$ and standard deviation $\sigma$ and take sufficiently large random samples from the population with replacement , then the distribution of the sample means will be approximately normally distributed.
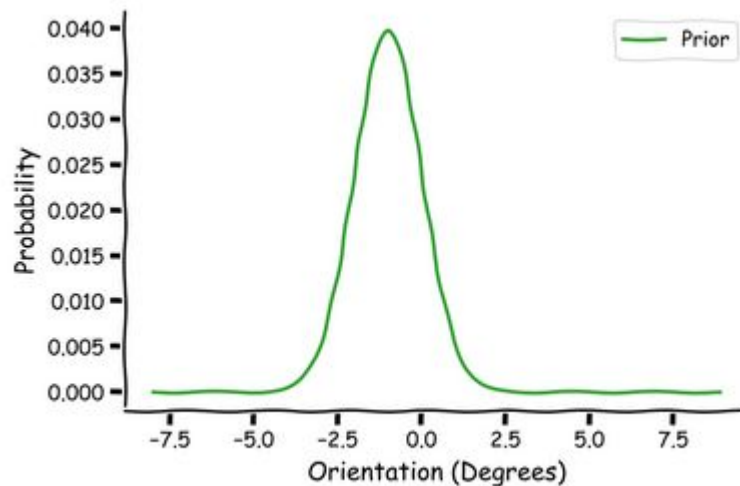


4

# Tutorial #1 Explanations

# Gaussian Distribution

Gaussians also have some mathematical/probabilistic (sum=1) properties that permit simple closed-form solutions to several important problems.

Gaussians have two parameters. The **mean** $\mu$, which sets the location of its center. Its "scale" or spread is controlled by its **standard deviation** $\sigma$ or its square, the **variance** $\sigma^2$.

$$\mathcal{N}(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-(x-\mu)^2}{2\sigma^2}\right)$$



6

Bayes Rule -

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{normalisation constant}}$$

Converting to guassians -

$$N\left(\mu_l, \sigma_l^2\right)$$

$$N\left(\mu_p, \sigma_p^2\right)$$

$$\text{posterior} \propto N\left(\mu_l, \sigma_l^2\right) \times N\left(\mu_p, \sigma_p^2\right)$$

(works for any distribution)

$$= N\left(\frac{\sigma_l^2 \mu_p + \sigma_p^2 \mu_l}{\sigma_p^2 + \sigma_l^2}, \frac{\sigma_p^2 \sigma_l^2}{\sigma_p^2 + \sigma_l^2}\right)$$

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

Closed form :

# Bayes theorem

TASK: estimate location of a noise-emitting object. To estimate its position, combine two sources of information:

1. new noisy auditory information (the likelihood)
2. prior visual expectations of where the stimulus is likely to come from (visual prior).

use Gaussian distributions to represent the auditory likelihood (in red), and a Gaussian visual prior (expectations - in blue). Using Bayes rule, combine them into a posterior distribution that summarizes the probability that the object is in each location.

Bayes' rule tells how to combine two sources of information: the prior (e.g., a noisy representation of our expectations about where the stimulus might come from) and the likelihood (e.g., a noisy representation of the stimulus position on a given trial), to obtain a posterior distribution taking into account both pieces of information.

W1D5_pod 031

prior : Noisy representation of our expectations about where stimulus might come from

likelihood : Noisy representation of stimulus position on a given trial

→ new noisy auditory information

visual prior : prior visual expectations of where stimulus is likely to come from!

combine them into posterior distribution that summarises the probability that object is in each location.

# Gaussian Parameter variance

Vary the parameters of Gaussians to see how changing the prior and likelihood affect the posterior.

# mu_auditory

# *mu_visual*

# sigma_visual

# Gaussians

The product of two Gaussian distributions, like our prior and likelihood, remains a Gaussian, regardless of the parameters. We can directly compute the parameters of that Gaussian from the means and variances of the prior and likelihood.

When does the prior have the strongest influence over the posterior? When is it the weakest?

Mu-values have the strongest influence over posterior. When combined with sigma (lower mu contributions) or when sigma is largest, prior seems to have a lesser influence over posterior.

# Conjugate-ness

Conjugate distributions or conjugate priors (for a particular likelihood) hold the following properties:

● The posterior has the same form (here, a normal distribution) as the prior, and

● There is simple, closed-form expression for its parameters.

Working with conjugate distributions is very convenient; otherwise, it is often necessary to use computationally-intensive numerical methods to combine the prior and likelihood.

product of gaussians → gaussian.

Compute parameters of that gaussian from means/
variances of that prior & likelihood!

$$\mu_p = \frac{\mu_{aud} \cdot \frac{1}{\sigma^2_{aud}} + \mu_{visual} \frac{1}{\sigma^2_{visual}}}{1/\sigma^2_{aud} + 1/\sigma^2_{visual}}$$

Sample output

Verifying conjugate properties.

① auditory likelihood constant

② compute posterior distribution
→ find mean

$$\int_x p(x)\,dx \quad \text{or} \quad \sum_x x \cdot p(x).$$

③ Compute analytical posterior mean
(from auditory & visual)

④ plot mean estimates

# Tutorial #1 Bonus Explanations

# Multimodal Priors

Gaussian prior: Stimulus is expected to come from a single location, though they might not know precisely where.

Multimodal priors: sound might come from one of two distinct locations.

We could model this using a Gaussian prior with a large $\sigma$ that covers both locations, but that would also make every point in between seem likely too. A better approach is to adjust the form of the prior so that it better matches the experiences/expectations by building a bimodal (2-peaked) prior out of Gaussians and examine the resulting posterior and its peaks.

Note: normalize the result so it is a proper probability distribution.

# Implementation and test!

Previous implementations do not help us answer questions like: What is the mean of our new prior? Is it a particularly likely location for the stimulus? Instead, we will use the posterior **mode** to summarize the distribution. The mode is the *location* of the most probable part of the distribution.

Observe what happens to the posterior as the likelihood gets closer to the different peaks of the prior.

Notice what happens to the posterior when the likelihood is exactly in between the two modes of the prior.



Example combination

Tutorial #2
Explanations

# Experimentation methodology

Study the effect of change as the distance between the visual stimulus (and the auditory stimulus increases/decreases

present only the auditory stimulus at varying locations and report where the source of the sound is located.using two pieces of information:

- The prior information about sound localization, learned during the trials before the curtain fell.
- Their noisy sensory estimates about where a particular sound originates.

The eventual goal is to predict the subjects' responses which implicitly requires building a prior that captures knowledge and expectations;

# Mixture of Gaussian Priors

single Gaussian prior -> could represent one of these possibilities.

A broad Gaussian with a large σ -> could represent sounds originating from nearly anywhere,

while a narrow Gaussian with μ near zero -> could represent sounds originating from the puppet.

Combine those into a mixture-of-Gaussians probability density function (PDF) that captures both possibilities

control Gaussian mixtures by summing them together with a 'mixing' or weight parameter pcommon, set to a value between 0 and 1

$$\text{Mixture} = P_c \times N(\mu_c, \sigma_c) + \left[ (1 - P_c) \times N(\mu_i, \tau_i) \right]$$

(common)

independent

probability that auditory stimulus shares a common source with the latent visual input

sounds independent of the first source

Note: $\quad P_c + P_i = 1 \Big\}$ law of total probability

gaussian_1 ≡ gaussian_common ($\mu = 0$, $\sigma = 0.5$)

gaussian_2 ≡ gaussian_independent ($\mu = 0$, $\sigma = 3$)

Combine ① & ② to make new prior by mixing

with mixing parameter $p_c = 0.75$

(peakier common cause gaussian
with 75% of the weight)

# Bayes theorem with Complex posteriors

We will compute the posterior by using Bayes Theorem to combine the mixture-of-gaussians prior and varying auditory Gaussian likelihood.

explore how a mixture-of-Gaussians prior and Gaussian likelihood interact

# Gaussian behavior Analysis

The mixture of Gaussian prior creates some interesting behaviour:
 1. We observe multiple modes (i.e. peaks) in our posterior
 (the common and independent causes).
 2. The mode of the posterior jumps between stimulus locations. These
 correspond to the participant switching from the independent to the common
 parts (i.e. causes) of the prior.

A similar discontinuity (ie. 'jump') in the posterior mode would happen in the
case of cue combination illusion with both auditory sources.
The same-source illusion breaks-down when the
voice stimulus is presented too far away from the visual input.

# Tutorial #3
# Explanations

# Gaussian behavior Analysis

compute all the necessary steps to perform model inversion (estimate the model parameters such as pcommon that generate data based on a mixture of Gaussian prior (common + independent priors) and a Gaussian likelihood similar to that of a participant).

for multiple possible stimulus inputs:

① create mixture of gaussian prior

② generate likelihood

③ estimate posterior as a $f(stimulus\_input)$

④ estimate participant response given posterior.

⑤ create distribution for input as $f(possible\ inputs)$

⑥ marginalisation

⑦ generate data using model (steps 1-4)

⑧ model inversion / model fitting using generated data to check for recovery of original parameters.

generative model {

model fitting / inversion! {

true visual stimulus

Stimulus

$x$

Brain

$\tilde{x}$

$p(\tilde{x}|x)$

Response

$\hat{x}$

$L(p(\tilde{x}|x))$
$+ p(x)$

$p(x|\tilde{x})$
(posterior)

prior

True stimulus
position

Noisy brain
encoding of
stimulus position

likelihood

Estimate of
stimulus position

**Known:**
True Position
presented in
experiment

**Unknown:**
Brain noisy
encoding of true
position

**Observable:**
Participant's
noisy estimation
response

ignored → (motor)

# Likelihood array

- *to consider all possible brain encodings*

multiple likelihood

$$f(x) = p(\tilde{x} \mid x)$$

# for each potential encoded stimulus $(\tilde{x})$

$\tilde{x}$

x

as a function of hypothesised true stimulus positions

encoded position

gaussian likelihood

$(\sigma_l = 1)$

1000

gaussian $\left(\overset{x}{points} \mid stim \mid \sigma\right)$

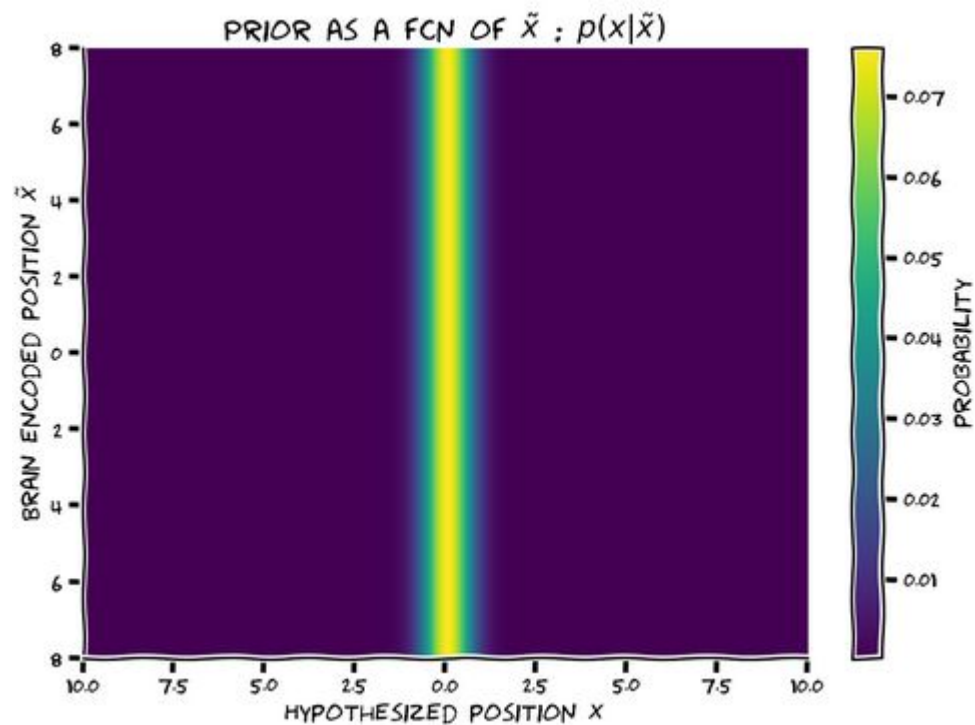each likelihood with different mean + different raw likelihood

# Causal mixture of gaussian prior

Create mixture of gaussian prior as $f$(brain encoded stimulus $\bar{x}$)

(prior doesn't change as $f(\tilde{x}) \Rightarrow$ identical $\forall$ row)

① create gaussian with $\mu = 0$, $\sigma = 0.5$.

② create gaussian; with $\mu = 0$, $\sigma = 10$.

③ combine ① & ② to make new prior } peakier gaussian has 95% of the weight

(mixing parameter $p_i = 0.05$)

# Bayes rule and posterior array

## HADAMARD PRODUCT

NumPy operations can often process an entire matrix in a single "vectorized" operation. This approach is often much faster and much easier to read than an element-by-element calculation.

$$\text{posterior}\ [i,:]\ \propto\ \text{likelihood}\ [i,:]\ \odot\ \text{Prior}\ [i,:]$$

for each encoded position $i$
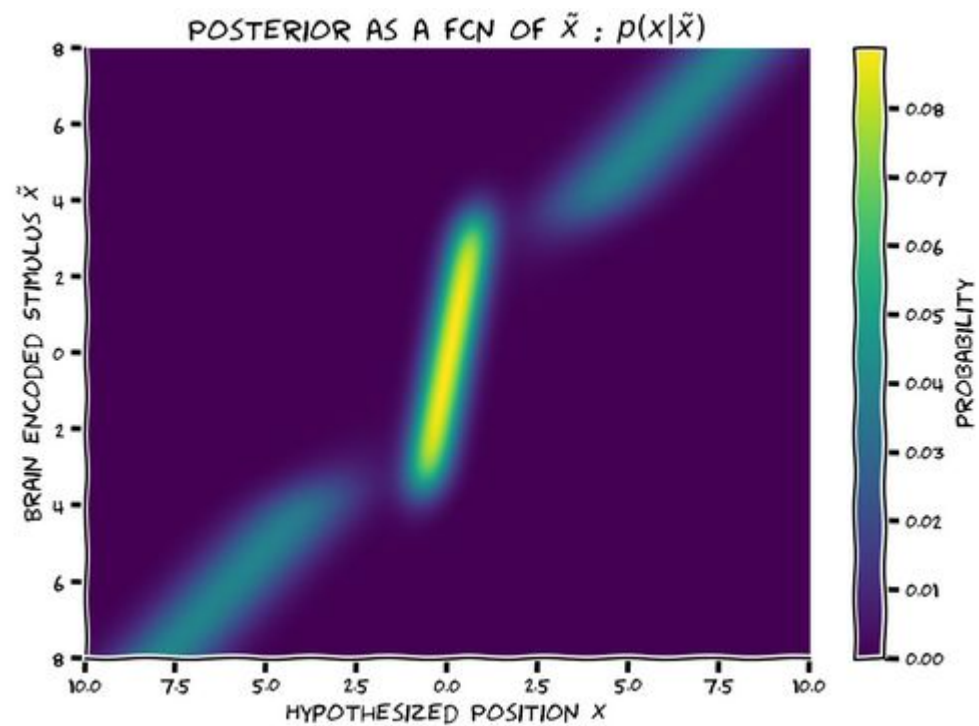
hadamard product
(element wise multiplication)

posterior distribution represents estimated stimulus position
$p(x|\tilde{x})$

represents posterior density

encodes decision for given encoding

mean of posterior = decision rule
(assumption)

estimate / response of sound location
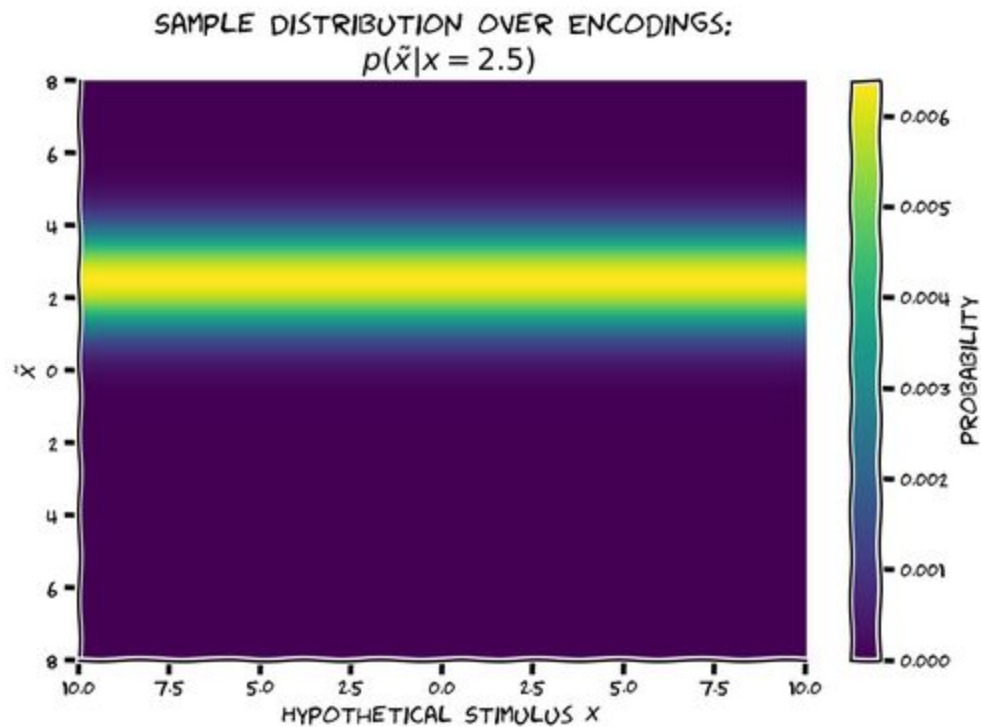$\hat{x}\ (f(\tilde{x}))$

POSTERIOR AS A FCN OF $\tilde{x}$ : $p(x|\tilde{x})$

Compute the binary decision array for each possible encoding

• calculate how likely each possible encoding is}
for given true stimulus

→ create gaussian centered around true presented stimulus
repeat across as $f$ (potentially encoded value $\check{x}$)   $(\sigma = 1)$

→ Column gaussian centered around true presented stimulus
repeat across as $f$ (hypothetical stimulus values $x$)

encodes distribution of brain encoded stimulus
& enable us to link true stimulus $x$ to
potential encoding $\check{x}$.

SAMPLE DISTRIBUTION OVER ENCODINGS:
$p(\tilde{x}|x = 2.5)$

true stimulus $x \longleftrightarrow$ potential encoding

$\downarrow$ calculate

distribution of encodings.

$\downarrow$

estimate

Marginalisation:

Marginalisation_Array = input array $\odot$ binary decision array (MA)

true presented stimulus

marginal $= \int_{\tilde{x}} MA$

all hypothetical values

MARGINALIZATION ARRAY: $p(\hat{x}|\tilde{x})$

# Data generation

Parameter recovery experiments are a powerful method for planning and debugging Bayesian analyses--if you cannot recover the given parameters, implementation has malfunctioned! Note that this value for pindependent is not quite the same as our prior, which used pindependent=0.05.



Participant behavior (')

mixing parameter

pind = 0.1

Model fitting algorithm -

① implement fria matlia
(recompute posterior/input/marginal) } depends on find

$\downarrow$
$p(\tilde{x}|x)$ { not likelihood as it
doesn't depend on find

② compute negative log likelihood ∀ trials
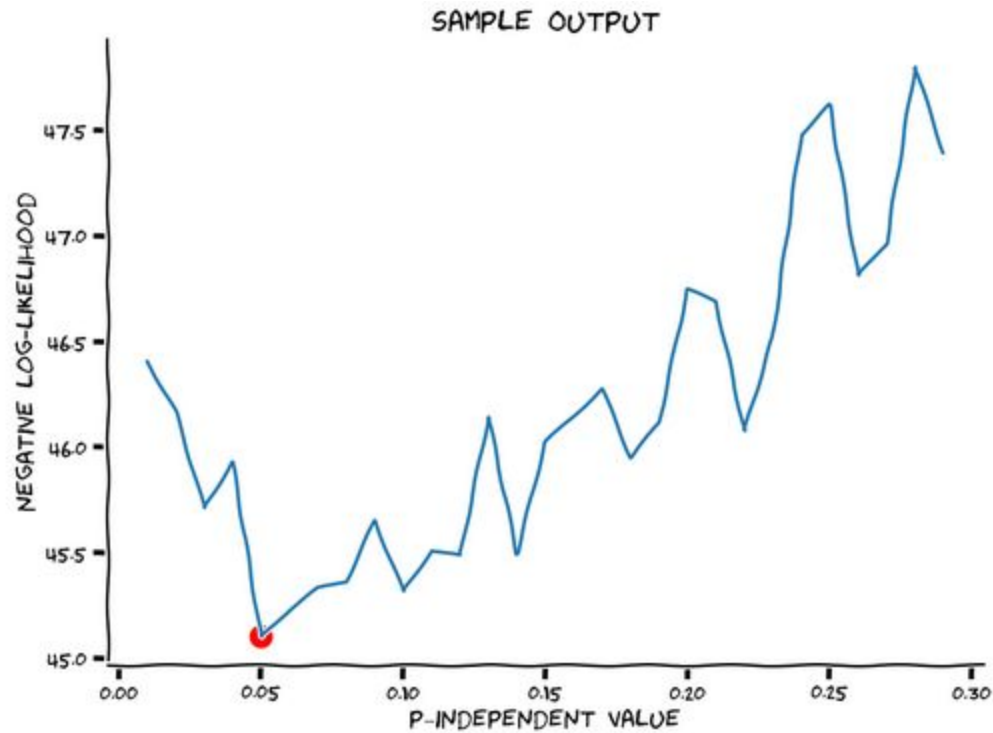find find that minimises -ve/max +ve
$(ll)$ $(ll)$

_assumption_ : trials are independent

trial i }

$$-Ll = -\sum_i log\, p\left(\hat{x}_i|x_i\right)$$

participant response.

$\longrightarrow$ presented stimulus

SAMPLE OUTPUT

# Conclusion

Importance of $p_{independent}$ = describes how much weight subjects assign to the same-cause vs. independent-cause origins of a sound

posterior: describes beliefs based on a combination of current evidence and prior experience.

Bayesian Analysis pipeline

develop model

↓

simulate data

↓

Bayes rule + marginalisation
(recover hidden parameters from the data)

# Tutorial #4 (Bonus) Explanations

# Agenda

1. Implement three commonly-used cost functions: mean-squared error, absolute error, and zero-one loss
2. Discover the concept of expected loss, and
3. Choose optimal locations on the posterior that minimize these cost functions and verify that these locations can be found analytically as well as empirically.
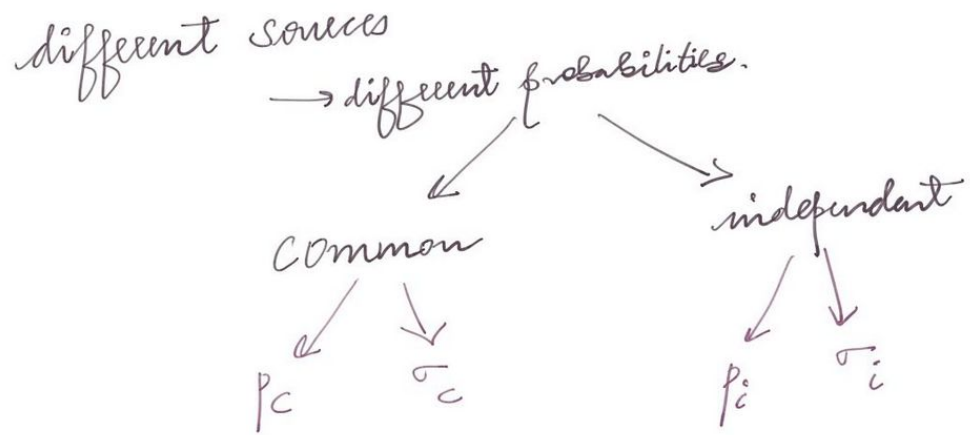
# Bayesian Decision theory

*combines the posterior with **cost functions** that allow us to quantify the potential impact of making a decision or choosing an action based on that posterior.*

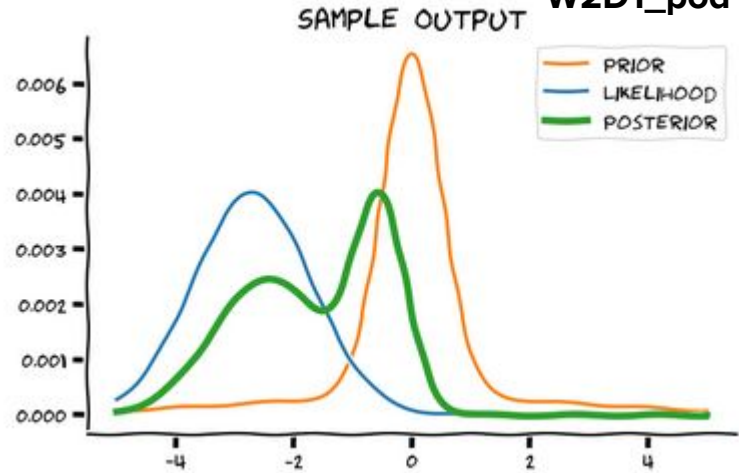probabilities ——— cost function → actions

Mean of posterior $p(x|\tilde{x})$ as proxy for response $\hat{x}$ for the participants

Using mean/median/mode of posterior distribution } as decision rule!

different sources

→ different probabilities.

common                  independant

$\rho_c$    $\sigma_c$         $\rho_i$    $\sigma_i$

$$prior = \begin{cases} N_c(0, 0.5) & \text{95\% weight} \\ N_i(0, 3) & \text{5\% weight} \end{cases}$$

$$likelihood = N(-2.7, 1)$$

SAMPLE OUTPUT

PRIOR
LIKELIHOOD
POSTERIOR

*Cost functions*



LOSS WHEN THE TRUE VALUE X=0

MEAN SQUARED ERROR
ABSOLUTE ERROR
ZERO-ONE LOSS

COST

PREDICTED VALUE $(\hat{x})$

Cost functions.

→ determines cost/penalty of estimating $\hat{x}$ when true/correct quantity $= x$

→ cost of error b/w true stimulus $x$ & estimate $\hat{x}$

① Mean squared error $= (x - \hat{x})^2$

② Absolute error $= |x - \hat{x}|$

③ Zero-one loss $= \begin{cases} 0 & \text{if } x = \hat{x} \\ 1 & \text{otherwise} \end{cases}$

# Expected loss function

A posterior distribution tells us about the confidence or credibility we assign to different choices. A cost function describes the penalty we incur when choosing an incorrect option. These concepts can be combined into an *expected loss* function
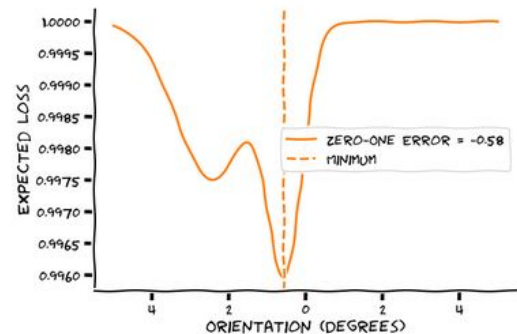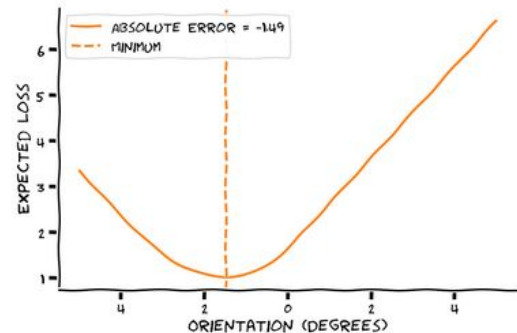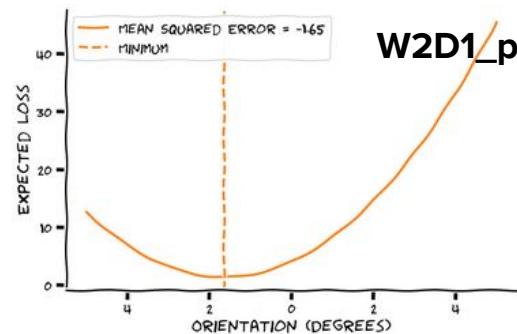
$$\text{Expected loss} \quad E\left[\text{loss} \mid \hat{x}\right] = \int L\left[\hat{x}, x\right] \odot p(x \mid \tilde{x}) \, dx$$

loss function

{ bimodal } posterior

hadamard product

(element wise multiplication)

# Expected loss outputs

minimum expected loss via brute-force: we searched over all

possible values of X and found the one that minimized

each of our loss functions.

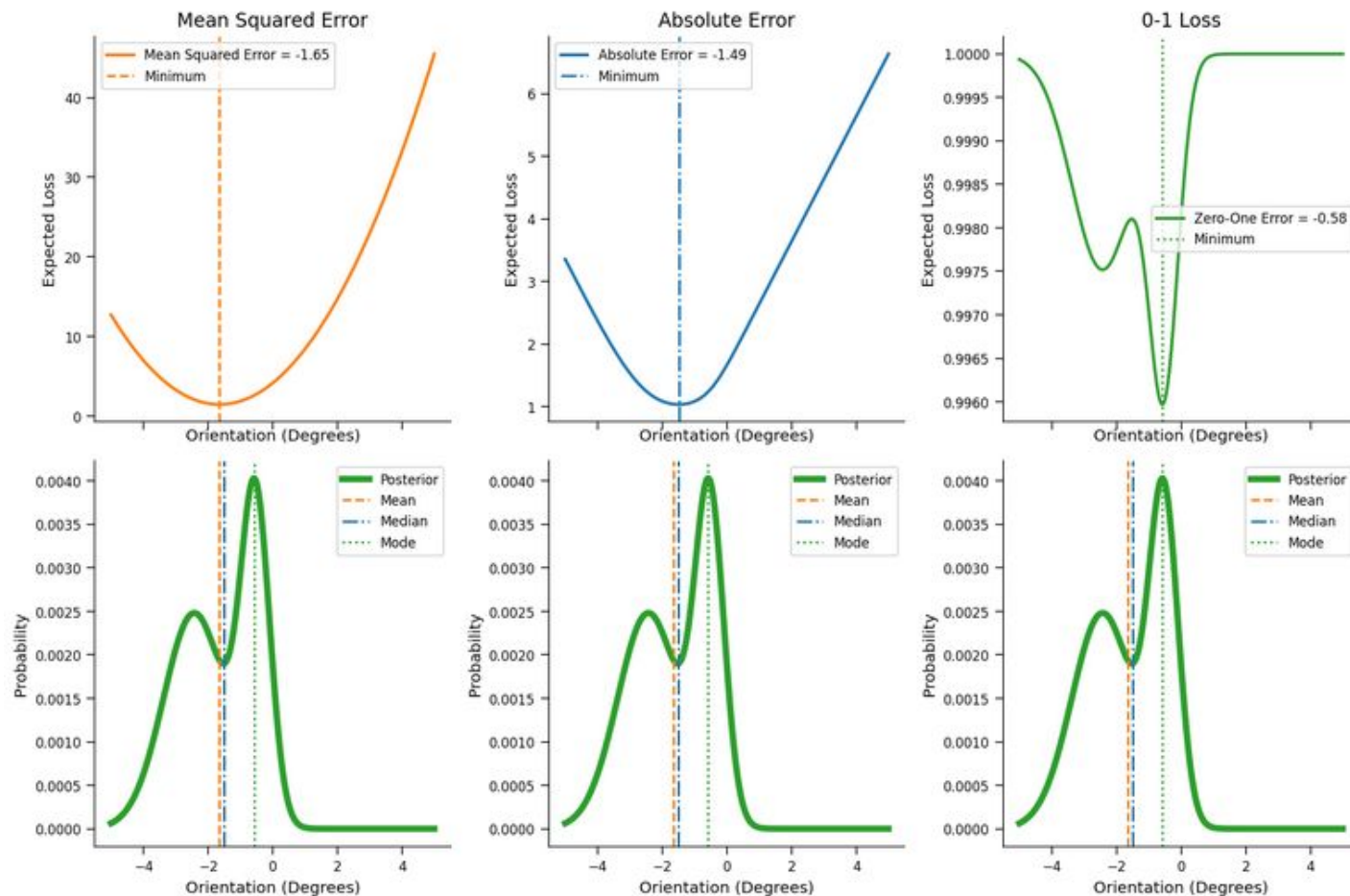This is feasible but can quickly become intractable.

# Summary

the three loss functions are minimized at specific points on the posterior, corresponding to the itss mean, median, and mode. To verify this property, replot the loss functions with the posterior on the same scale beneath. The mean, median, and mode are marked on the posterior.

We used expected loss to quantify the results of making a decision, and showed that optimizing under different cost functions led us to choose different locations on the posterior. Finally, we found that these optimal locations can be identified analytically, sparing us from a brute-force search.

OBSERVATIONS: The mean minimizes the mean-squared error.

Absolute error is minimized by the median,

while zero-one loss is minimized at the posterior's mode.

# Food for thought

- *Suppose your professor offered to grade your work with a zero-one loss or mean square error.*
  - *When might you choose each?*
  - *Which would be easier to learn from?*

- *All of the loss functions we considered are symmetrical. Are there situations where an asymmetrical loss function might make sense? How about a negative one?*

# Question #1 Answer

Firstly, the mean squared error is close to the variance, however you average the value of variance out by the number of the observations. In a way, it is a mean//average//expected value of the variance//dispersion of the data values.

The sum of your losses would no longer represent accuracy in this case, but rather the total "cost" of misclassification. The 0-1 loss function is unique in its equivalence to accuracy, since all you care about is whether you got it right or not, and not how the errors are made.

So, prefer mean squared errors while it's probably easier to learn from 0-1 loss as it's more straightforward.

# Question #2 Answer

Not all loss is the same. So, we weight different losses in the loss functions giving rise to an asymmetrical loss metric.

It is the case that we often use loss functions that become equal to zero when the fit of the model to the training data is perfect, but the optimization algorithms don't care about this, and they drive the loss function to algebraically more negative values, and not towards zero.

# References

http://www.lenstrnad.com/blog/2018/09/asymmetric_loss_function

https://discuss.pytorch.org/t/what-happens-when-loss-are-negative/47883/3

https://stats.stackexchange.com/questions/284028/0-1-loss-function-explanation/284062