

Decision problems

September 4, 2019

- 1 Beliefs and probabilities
 - Probability and Bayesian inference
- 2 Hierarchies of decision making problems
- 3 Formalising Classification problems
- 4 Classification with stochastic gradient descent

Uncertainty

- We cannot perfectly predict the future.
- We cannot know for sure what happened in the past.
- How can we quantify this uncertainty?
- Probabilities!

Axioms of probability

For any probability measure^a P on (Ω, Σ) ,

^a Σ is the set of possible events, with $A \in \Sigma$ always $A \subset \Omega$. Technically Σ is a σ -algebra

Uncertainty

- We cannot perfectly predict the future.
- We cannot know for sure what happened in the past.
- How can we quantify this uncertainty?
- Probabilities!

Axioms of probability

For any probability measure^a P on (Ω, Σ) ,

- 1 The probability of the certain event is $P(\Omega) = 1$

^a Σ is the set of possible events, with $A \in \Sigma$ always $A \subset \Omega$. Technically Σ is a σ -algebra

Uncertainty

- We cannot perfectly predict the future.
- We cannot know for sure what happened in the past.
- How can we quantify this uncertainty?
- Probabilities!

Axioms of probability

For any probability measure^a P on (Ω, Σ) ,

- 1 The probability of the certain event is $P(\Omega) = 1$
- 2 The probability of the impossible event is $P(\emptyset) = 0$

^a Σ is the set of possible events, with $A \in \Sigma$ always $A \subset \Omega$. Technically Σ is a σ -algebra

Uncertainty

- We cannot perfectly predict the future.
- We cannot know for sure what happened in the past.
- How can we quantify this uncertainty?
- Probabilities!

Axioms of probability

For any probability measure^a P on (Ω, Σ) ,

- 1 The probability of the certain event is $P(\Omega) = 1$
- 2 The probability of the impossible event is $P(\emptyset) = 0$
- 3 The probability of any event $A \in \Sigma$ is $0 \leq P(A) \leq 1$.

^a Σ is the set of possible events, with $A \in \Sigma$ always $A \subset \Omega$. Technically Σ is a σ -algebra

Uncertainty

- We cannot perfectly predict the future.
- We cannot know for sure what happened in the past.
- How can we quantify this uncertainty?
- Probabilities!

Axioms of probability

For any probability measure^a P on (Ω, Σ) ,

- 1 The probability of the certain event is $P(\Omega) = 1$
- 2 The probability of the impossible event is $P(\emptyset) = 0$
- 3 The probability of any event $A \in \Sigma$ is $0 \leq P(A) \leq 1$.
- 4 If A, B are disjoint, i.e. $A \cap B = \emptyset$, meaning that they cannot happen at the same time, then

$$P(A \cup B) = P(A) + P(B)$$

^a Σ is the set of possible events, with $A \in \Sigma$ always $A \subset \Omega$. Technically Σ is a σ -algebra

Definition 1 (Conditional probability)

The probability of A happening if we know that B has happened is defined to be:

$$P(A | B) \triangleq \frac{P(A \cap B)}{P(B)}.$$

Conditional probabilities obey the same rules as probabilities.

Bayes's theorem

For $P(A_1 \cup A_2) = 1$, $A_1 \cap A_2 = \emptyset$,

$$P(A_i | B)$$

Definition 1 (Conditional probability)

The probability of A happening if we know that B has happened is defined to be:

$$P(A | B) \triangleq \frac{P(A \cap B)}{P(B)}.$$

Bayes's theorem

For $P(A_1 \cup A_2) = 1$, $A_1 \cap A_2 = \emptyset$,

$$P(A_i | B) = \frac{P(B | A_i)P(A_i)}{P(B)}$$

Definition 1 (Conditional probability)

The probability of A happening if we know that B has happened is defined to be:

$$P(A | B) \triangleq \frac{P(A \cap B)}{P(B)}.$$

Bayes's theorem

For $P(A_1 \cup A_2) = 1$, $A_1 \cap A_2 = \emptyset$,

$$P(A_i | B) = \frac{P(B | A_i)P(A_i)}{P(B)} = \frac{P(B | A_i)P(A_i)}{P(B | A_1)P(A_1) + P(B | A_2)P(A_2)}$$

Definition 1 (Conditional probability)

The probability of A happening if we know that B has happened is defined to be:

$$P(A | B) \triangleq \frac{P(A \cap B)}{P(B)}.$$

Bayes's theorem

For $P(A_1 \cup A_2) = 1$, $A_1 \cap A_2 = \emptyset$,

$$P(A_i | B) = \frac{P(B | A_i)P(A_i)}{P(B)} = \frac{P(B | A_i)P(A_i)}{P(B | A_1)P(A_1) + P(B | A_2)P(A_2)}$$

Example 2 (probability of rain)

What is the probability of rain given a forecast x_1 or x_2 ?

$$\begin{array}{l|l} \omega_1: \text{rain} & P(\omega_1) = 80\% \\ \omega_2: \text{dry} & P(\omega_2) = 20\% \end{array}$$

Table: Prior probability of rain tomorrow

Definition 1 (Conditional probability)

The probability of A happening if we know that B has happened is defined to be:

$$P(A | B) \triangleq \frac{P(A \cap B)}{P(B)}.$$

Bayes's theorem

For $P(A_1 \cup A_2) = 1$, $A_1 \cap A_2 = \emptyset$,

$$P(A_i | B) = \frac{P(B | A_i)P(A_i)}{P(B)} = \frac{P(B | A_i)P(A_i)}{P(B | A_1)P(A_1) + P(B | A_2)P(A_2)}$$

Example 2 (probability of rain)

What is the probability of rain given a forecast x_1 or x_2 ?

$$\begin{array}{l|l} \omega_1: \text{rain} & P(\omega_1) = 80\% \\ \omega_2: \text{dry} & P(\omega_2) = 20\% \end{array}$$

$$\begin{array}{l|l} x_1: \text{rain} & P(x_1 | \omega_1) = 90\% \\ x_2: \text{dry} & P(x_2 | \omega_2) = 50\% \end{array}$$

Table: Prior probability of rain tomorrow

Table: Probability the forecast is correct

Definition 1 (Conditional probability)

The probability of A happening if we know that B has happened is defined to be:

$$P(A | B) \triangleq \frac{P(A \cap B)}{P(B)}.$$

Bayes's theorem

For $P(A_1 \cup A_2) = 1$, $A_1 \cap A_2 = \emptyset$,

$$P(A_i | B) = \frac{P(B | A_i)P(A_i)}{P(B)} = \frac{P(B | A_i)P(A_i)}{P(B | A_1)P(A_1) + P(B | A_2)P(A_2)}$$

Example 2 (probability of rain)

What is the probability of rain given a forecast x_1 or x_2 ?

$$\begin{array}{l|l} \omega_1: \text{rain} & P(\omega_1) = 80\% \\ \omega_2: \text{dry} & P(\omega_2) = 20\% \end{array}$$

$$\begin{array}{l|l} x_1: \text{rain} & P(x_1 | \omega_1) = 90\% \\ x_2: \text{dry} & P(x_2 | \omega_2) = 50\% \end{array}$$

$$\begin{array}{l} P(\omega_1 | x_1) = 87.8\% \\ P(\omega_1 | x_2) = 44.4\% \end{array}$$

Table: Prior probability of rain tomorrow

Table: Probability the forecast is correct

Table: Probability that it will rain given the forecast

Classification in terms of conditional probabilities

- Features $x_t \in \mathcal{X}$.
- Class label $y_t \in \mathcal{Y}$.
- Probability model $P_\mu(x_t | y_t)$.
- Prior class probability $P_\mu(y_t = c)$.

$$P_\mu(y_t = c | x_t) = \frac{P_\mu(x_t | y_t = c)P_\mu(y_t = c)}{\sum_{c' \in \mathcal{Y}} P_\mu(x_t | y_t = c')P_\mu(y_t = c')}$$

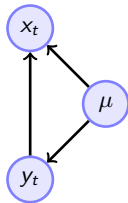
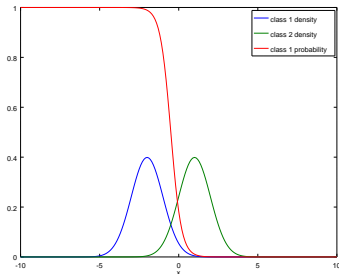


Figure: A generative classification model. μ identifies the model (parameter). x_t are the features and y_t the class label of the t -th example.

Classification in terms of conditional probabilities



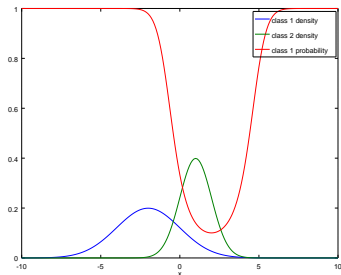
(a) Equal prior and variance

Figure: The effect of changing variance and prior when we assume a normal distribution.

Example 3 (Normal distribution)

A simple example is when x_t is normally distributed in a manner that depends on the class. Figure 2 shows the distribution of x_t for two different classes, with means of -1 and $+1$ respectively, for three different case. In the first case, both classes have variance of 1, and we assume the same prior probability for both

Classification in terms of conditional probabilities



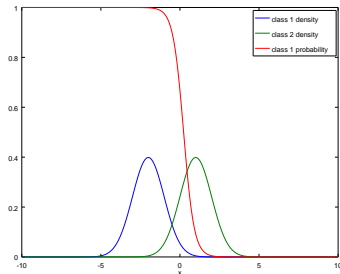
(a) Unequal variance

Figure: The effect of changing variance and prior when we assume a normal distribution.

Example 3 (Normal distribution)

A simple example is when x_t is normally distributed in a manner that depends on the class. Figure 2 shows the distribution of x_t for two different classes, with means of -1 and $+1$ respectively, for three different cases. In the first case, both classes have variance of 1, and we assume the same prior probability for both

Classification in terms of conditional probabilities



(a) Unequal prior

Figure: The effect of changing variance and prior when we assume a normal distribution.

Example 3 (Normal distribution)

A simple example is when x_t is normally distributed in a manner that depends on the class. Figure 2 shows the distribution of x_t for two different classes, with means of -1 and $+1$ respectively, for three different cases. In the first case, both classes have variance of 1, and we assume the same prior probability for both

Classification in terms of conditional probabilities

Figure: The effect of changing variance and prior when we assume a normal distribution.

Example 3 (Normal distribution)

A simple example is when x_t is normally distributed in a manner that depends on the class. Figure 2 shows the distribution of x_t for two different classes, with means of -1 and $+1$ respectively, for three different cases. In the first case, both classes have variance of 1, and we assume the same prior probability for both

$$x_t \mid y_t = 0 \sim \mathcal{N}(-1, 1), \quad x_t \mid y_t = 1 \sim \mathcal{N}(1, 1)$$

$$x_t \mid y_t = 0 \sim \mathcal{N}(-1, 1), \quad x_t \mid y_t = 1 \sim \mathcal{N}(1, 1)$$

But how can we get a probability model in the first place?

Subjective probability

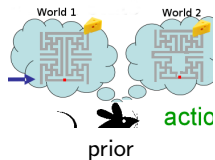
Subjective probability measure ξ

- If we think event A is more likely than B , then $\xi(A) > \xi(B)$.
- Usual rules of probability apply:
 - 1 $\xi(A) \in [0, 1]$.
 - 2 $\xi(\emptyset) = 0$.
 - 3 If $A \cap B = \emptyset$, then $\xi(A \cup B) = \xi(A) + \xi(B)$.

Bayesian inference illustration

Use a subjective belief $\xi(\mu)$ on \mathcal{M}

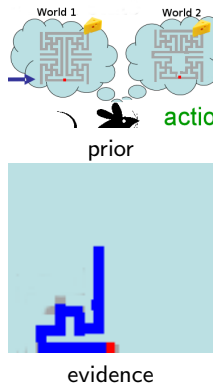
- **Prior** belief $\xi(\mu)$ represents our initial uncertainty.



Bayesian inference illustration

Use a subjective belief $\xi(\mu)$ on \mathcal{M}

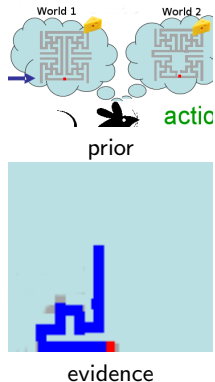
- **Prior** belief $\xi(\mu)$ represents our initial uncertainty.
- We **observe history** h .



Bayesian inference illustration

Use a subjective belief $\xi(\mu)$ on \mathcal{M}

- **Prior** belief $\xi(\mu)$ represents our initial uncertainty.
- We **observe history** h .
- Each possible μ assigns a **probability** $P_\mu(h)$ to h .

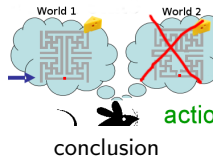
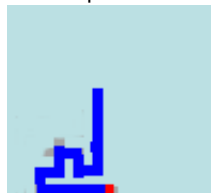
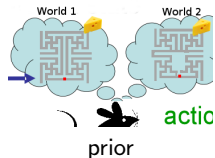


Bayesian inference illustration

Use a subjective belief $\xi(\mu)$ on \mathcal{M}

- **Prior** belief $\xi(\mu)$ represents our initial uncertainty.
- We **observe history** h .
- Each possible μ assigns a **probability** $P_\mu(h)$ to h .
- We can use this to **update** our belief via Bayes' theorem to obtain the **posterior** belief:

$$\xi(\mu | h) \propto P_\mu(h)\xi(\mu) \quad (\text{conclusion} = \text{evidence} \times \text{prior})$$



Some examples

Example 4

John claims to be a medium. He throws a coin n times and predicts its value always correctly. Should we believe that he is a medium?

- μ_1 : John is a medium.
- μ_0 : John is not a medium.

The answer depends on what we **expect** a medium to be able to do, and how likely we thought he'd be a medium in the first place.

Bayesian inference

- mutually exclusive models $\mathcal{M} = \{\mu_1, \dots, \mu_k\}$.

Bayesian inference

- mutually exclusive models $\mathcal{M} = \{\mu_1, \dots, \mu_k\}$.
- Probability model for any data x : $P_\mu(x) \equiv \mathbb{P}(x \mid \mu)$.

Bayesian inference

- mutually exclusive models $\mathcal{M} = \{\mu_1, \dots, \mu_k\}$.
- Probability model for any data x : $P_\mu(x) \equiv \mathbb{P}(x \mid \mu)$.
- For each model, we have a prior probability $\xi(\mu)$ that it is correct.

Bayesian inference

- mutually exclusive models $\mathcal{M} = \{\mu_1, \dots, \mu_k\}$.
- Probability model for any data x : $P_\mu(x) \equiv \mathbb{P}(x | \mu)$.
- For each model, we have a prior probability $\xi(\mu)$ that it is correct.
- Posterior probability

$$\xi(\mu | x) = \frac{\mathbb{P}(x | \mu)\xi(\mu)}{\sum_{\mu' \in \mathcal{M}} \mathbb{P}(x | \mu')\xi(\mu')} = \frac{P_\mu(x)\xi(\mu)}{\sum_{\mu' \in \mathcal{M}} P_{\mu'}(x)\xi(\mu')}.$$

Bayesian inference

- mutually exclusive models $\mathcal{M} = \{\mu_1, \dots, \mu_k\}$.
- Probability model for any data x : $P_\mu(x) \equiv \mathbb{P}(x | \mu)$.
- For each model, we have a prior probability $\xi(\mu)$ that it is correct.
- Posterior probability

$$\xi(\mu | x) = \frac{\mathbb{P}(x | \mu)\xi(\mu)}{\sum_{\mu' \in \mathcal{M}} \mathbb{P}(x | \mu')\xi(\mu')} = \frac{P_\mu(x)\xi(\mu)}{\sum_{\mu' \in \mathcal{M}} P_{\mu'}(x)\xi(\mu')}.$$

Interpretation

- \mathcal{M} : Set of all possible models that could describe the data.
- $P_\mu(x)$: Probability of x under model μ .
- Alternative notation $\mathbb{P}(x | \mu)$: Probability of x given that model μ is correct.
- $\xi(\mu)$: Our belief, before seeing the data, that μ is correct.
- $\xi(\mu | x)$: Our belief, after seeing the data, that μ is correct.

Exercise 1 (Continued example for medium)

$$P_{\mu}(x) = \prod_{t=1}^n P_{\mu}(x_t). \quad (\text{independence property})$$

$$P_{\mu_1}(x_t = 1) = 1, \quad P_{\mu_1}(x_t = 0) = 0. \quad (\text{true medium model})$$

$$P_{\mu_0}(x_t = 1) = 1/2, \quad P_{\mu_0}(x_t = 0) = 1/2. \quad (\text{non-medium model})$$

Throw a coin 4 times, and have a classmate make a prediction. What your belief that your classmate is a medium? Is the prior you used reasonable?

Exercise 1 (Continued example for medium)

$$P_{\mu}(x) = \prod_{t=1}^n P_{\mu}(x_t). \quad (\text{independence property})$$

$$P_{\mu_1}(x_t = 1) = 1, \quad P_{\mu_1}(x_t = 0) = 0. \quad (\text{true medium model})$$

$$P_{\mu_0}(x_t = 1) = 1/2, \quad P_{\mu_0}(x_t = 0) = 1/2. \quad (\text{non-medium model})$$

Throw a coin 4 times, and have a classmate make a prediction. What your belief that your classmate is a medium? Is the prior you used reasonable?

Exercise 1 (Continued example for medium)

$$P_{\mu}(x) = \prod_{t=1}^n P_{\mu}(x_t). \quad (\text{independence property})$$

$$P_{\mu_1}(x_t = 1) = 1, \quad P_{\mu_1}(x_t = 0) = 0. \quad (\text{true medium model})$$

$$P_{\mu_0}(x_t = 1) = 1/2, \quad P_{\mu_0}(x_t = 0) = 1/2. \quad (\text{non-medium model})$$

$$\xi(\mu_0) = 1/2, \quad \xi(\mu_1) = 1/2. \quad (\text{prior belief})$$

Throw a coin 4 times, and have a classmate make a prediction. What your belief that your classmate is a medium? Is the prior you used reasonable?

Exercise 1 (Continued example for medium)

$$P_{\mu}(x) = \prod_{t=1}^n P_{\mu}(x_t). \quad (\text{independence property})$$

$$P_{\mu_1}(x_t = 1) = 1, \quad P_{\mu_1}(x_t = 0) = 0. \quad (\text{true medium model})$$

$$P_{\mu_0}(x_t = 1) = 1/2, \quad P_{\mu_0}(x_t = 0) = 1/2. \quad (\text{non-medium model})$$

$$\xi(\mu_0) = 1/2, \quad \xi(\mu_1) = 1/2. \quad (\text{prior belief})$$

$$\xi(\mu_1 | x) = \frac{P_{\mu_1}(x)\xi(\mu_1)}{\mathbb{P}_{\xi}(x)} \quad (\text{posterior belief})$$

$$\mathbb{P}_{\xi}(x) \triangleq P_{\mu_1}(x)\xi(\mu_1) + P_{\mu_0}(x)\xi(\mu_0). \quad (\text{marginal distribution})$$

Throw a coin 4 times, and have a classmate make a prediction. What your belief that your classmate is a medium? Is the prior you used reasonable?

Sequential update of beliefs

	M	T	W	T	F	S	S
CNN	0.5	0.6	0.7	0.9	0.5	0.3	0.1
SMHI	0.3	0.7	0.8	0.9	0.5	0.2	0.1
YR	0.6	0.9	0.8	0.5	0.4	0.1	0.1
Rain?	Y	Y	Y	N	Y	N	N

Table: Predictions by three different entities for the probability of rain on a particular day, along with whether or not it actually rained.

Exercise 2

- n meteorological stations $\{\mu_i \mid i = 1, \dots, n\}$
- The i -th station predicts rain $P_{\mu_i}(x_t \mid x_1, \dots, x_{t-1})$.
- Let $\xi_t(\mu)$ be our belief at time t . Derive the next-step belief $\xi_{t+1}(\mu) \triangleq \xi_t(\mu|y_t)$ in terms of the current belief ξ_t .
- Write a python function that computes this posterior

Sequential update of beliefs

	M	T	W	T	F	S	S
CNN	0.5	0.6	0.7	0.9	0.5	0.3	0.1
SMHI	0.3	0.7	0.8	0.9	0.5	0.2	0.1
YR	0.6	0.9	0.8	0.5	0.4	0.1	0.1
Rain?	Y	Y	Y	N	Y	N	N

Table: Predictions by three different entities for the probability of rain on a particular day, along with whether or not it actually rained.

Exercise 2

- n meteorological stations $\{\mu_i \mid i = 1, \dots, n\}$
- The i -th station predicts rain $P_{\mu_i}(x_t \mid x_1, \dots, x_{t-1})$.
- Let $\xi_t(\mu)$ be our belief at time t . Derive the next-step belief $\xi_{t+1}(\mu) \triangleq \xi_t(\mu|y_t)$ in terms of the current belief ξ_t .
- Write a python function that computes this posterior

$$\xi_{t+1}(\mu) \triangleq \xi_t(\mu|x_t) = \frac{P_{\mu}(x_t \mid x_1, \dots, x_{t-1})\xi_t(\mu)}{\sum_{\mu'} P_{\mu'}(x_t \mid x_1, \dots, x_{t-1})\xi_t(\mu')}$$

Bayesian inference for Bernoulli distributions

Estimating a coin's bias

A fair coin comes heads 50% of the time. We want to test an unknown coin, which we think may not be completely fair.

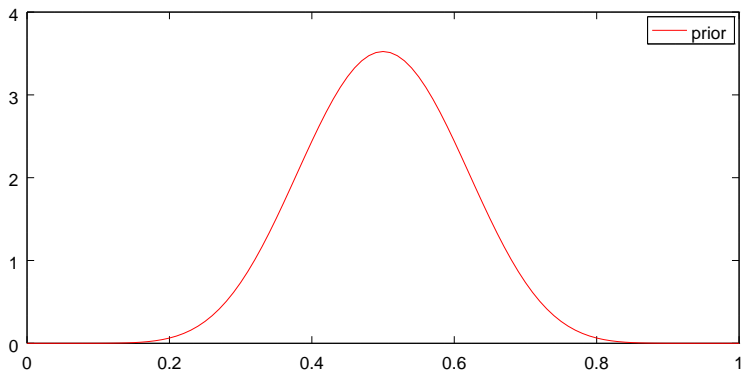


Figure: Prior belief ξ about the coin bias θ .

Bayesian inference for Bernoulli distributions

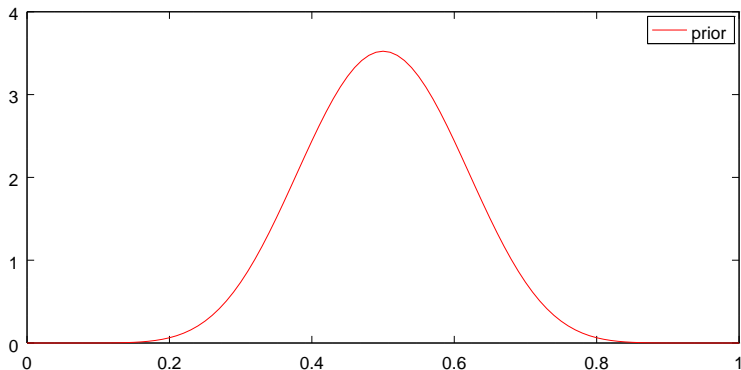


Figure: Prior belief ξ about the coin bias θ .

For a sequence of throws $x_t \in \{0, 1\}$,

$$P_{\theta}(x) \propto \prod_t \theta^{x_t} (1 - \theta)^{1 - x_t} = \theta^{\#\text{Heads}} (1 - \theta)^{\#\text{Tails}}$$

Bayesian inference for Bernoulli distributions

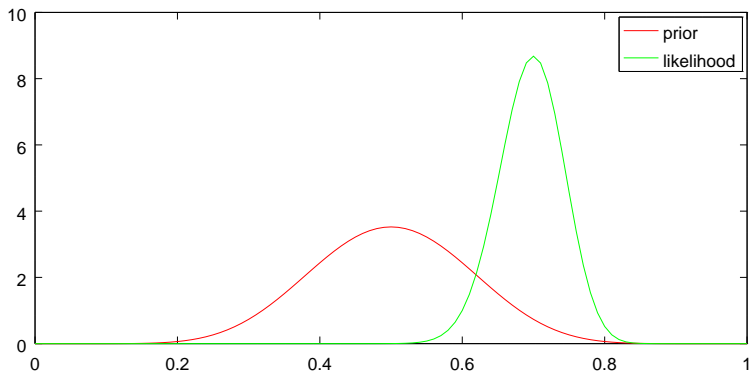


Figure: Prior belief ξ about the coin bias θ and likelihood of θ for the data.

Say we throw the coin 100 times and obtain 70 heads. Then we plot the **likelihood** $P_{\theta}(x)$ of different models.

Bayesian inference for Bernoulli distributions

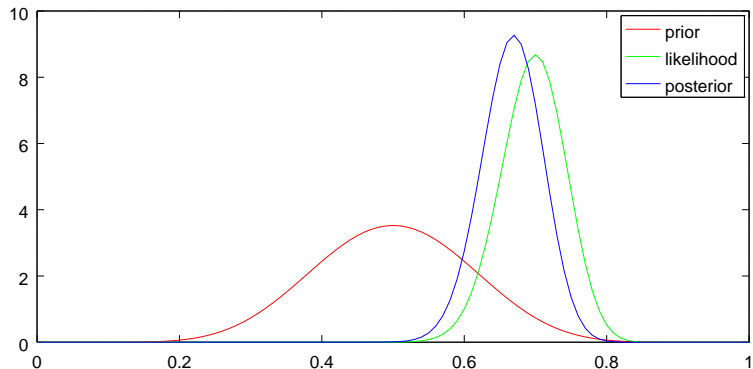


Figure: Prior belief $\xi(\theta)$ about the coin bias θ , likelihood of θ for the data, and posterior belief $\xi(\theta | x)$

From these, we calculate a **posterior** distribution over the correct models. This represents our conclusion given our prior and the data.

Learning outcomes

Understanding

- The axioms of probability, marginals and conditional distributions.
- The philosophical underpinnings of Bayesianism.
- The simple conjugate model for Bernoulli distributions.

Skills

- Be able to calculate with probabilities using the marginal and conditional definitions and Bayes rule.
- Being able to implement a simple Bayesian inference algorithm in Python.

Reflection

- How useful is the Bayesian representation of uncertainty?
- How restrictive is the need to select a prior distribution?
- Can you think of another way to explicitly represent uncertainty in a way that can incorporate new evidence?

- 1 Beliefs and probabilities
- 2 Hierarchies of decision making problems**
 - Simple decision problems
 - Decision rules
- 3 Formalising Classification problems
- 4 Classification with stochastic gradient descent

Preferences

Example 5

Food

- A McDonald's cheeseburger
- B Surstromming
- C Oatmeal

Money

- A 10,000,000 SEK
- B 10,000,000 USD
- C 10,000,000 BTC

Entertainment

- A Ticket to Liseberg
- B Ticket to Rebstar
- C Ticket to Nutcracker

Rewards and utilities

- Each choice is called a **reward** $r \in \mathcal{R}$.
- There is a **utility function** $U : \mathcal{R} \rightarrow \mathbb{R}$, assigning values to reward.
- We (weakly) prefer A to B iff $U(A) \geq U(B)$.

Exercise 3

From your individual preferences, derive a **common utility function** that reflects everybody's preferences in the class for each of the three examples. Is there a simple algorithm for deciding this? Would you consider the outcome fair?

Preferences among random outcomes

Example 6

Would you rather . . .

- A Have 100 EUR now?
- B Flip a coin, and get 200 EUR if it comes heads?

Risk and monetary rewards

Preferences among random outcomes

Example 6

Would you rather ...

- A Have 100 EUR now?
- B Flip a coin, and get 200 EUR if it comes heads?

The expected utility hypothesis

Rational decision makers prefer choice A to B if

$$\mathbb{E}(U|A) \geq \mathbb{E}(U|B),$$

where the expected utility is

$$\mathbb{E}(U|A) = \sum_r U(r) \mathbb{P}(r|A).$$

In the above example, $r \in \{0, 100, 200\}$ and $U(r)$ is increasing, and the coin is fair.

Risk and monetary rewards

Preferences among random outcomes

Example 6

Would you rather ...

- A Have 100 EUR now?
- B Flip a coin, and get 200 EUR if it comes heads?

The expected utility hypothesis

Rational decision makers prefer choice A to B if

$$\mathbb{E}(U|A) \geq \mathbb{E}(U|B),$$

where the expected utility is

$$\mathbb{E}(U|A) = \sum_r U(r) \mathbb{P}(r|A).$$

In the above example, $r \in \{0, 100, 200\}$ and $U(r)$ is increasing, and the coin is fair.

Risk and monetary rewards

- If U is convex, we are risk-seeking.

Preferences among random outcomes

Example 6

Would you rather ...

- A Have 100 EUR now?
- B Flip a coin, and get 200 EUR if it comes heads?

The expected utility hypothesis

Rational decision makers prefer choice A to B if

$$\mathbb{E}(U|A) \geq \mathbb{E}(U|B),$$

where the expected utility is

$$\mathbb{E}(U|A) = \sum_r U(r) \mathbb{P}(r|A).$$

In the above example, $r \in \{0, 100, 200\}$ and $U(r)$ is increasing, and the coin is fair.

Risk and monetary rewards

- If U is convex, we are risk-seeking.
- If U is concave, we are risk-averse.

Preferences among random outcomes

Example 6

Would you rather ...

- A Have 100 EUR now?
- B Flip a coin, and get 200 EUR if it comes heads?

The expected utility hypothesis

Rational decision makers prefer choice A to B if

$$\mathbb{E}(U|A) \geq \mathbb{E}(U|B),$$

where the expected utility is

$$\mathbb{E}(U|A) = \sum_r U(r) \mathbb{P}(r|A).$$

In the above example, $r \in \{0, 100, 200\}$ and $U(r)$ is increasing, and the coin is fair.

Risk and monetary rewards

- If U is convex, we are risk-seeking.
- If U is linear, we are risk neutral.
- If U is concave, we are risk-averse.

Uncertain rewards

- Decisions $a \in \mathcal{A}$
- Each choice is called a **reward** $r \in \mathcal{R}$.
- There is a **utility function** $U : \mathcal{R} \rightarrow \mathbb{R}$, assigning values to reward.
- We (weakly) prefer A to B iff $U(A) \geq U(B)$.

Example 7

You are going to work, and it might rain.
What do you do?

- a_1 : Take the umbrella.
- a_2 : Risk it!
- ω_1 : rain
- ω_2 : dry

$\rho(\omega, a)$	a_1	a_2
ω_1	dry, carrying umbrella	wet
ω_2	dry, carrying umbrella	dry
$U[\rho(\omega, a)]$	a_1	a_2
ω_1	0	-10
ω_2	0	1

Table: Rewards and utilities.

Uncertain rewards

- Decisions $a \in \mathcal{A}$
- Each choice is called a **reward** $r \in \mathcal{R}$.
- There is a **utility function** $U : \mathcal{R} \rightarrow \mathbb{R}$, assigning values to reward.
- We (weakly) prefer A to B iff $U(A) \geq U(B)$.

Example 7

You are going to work, and it might rain.
What do you do?

- a_1 : Take the umbrella.
- a_2 : Risk it!
- ω_1 : rain
- ω_2 : dry

$\rho(\omega, a)$	a_1	a_2
ω_1	dry, carrying umbrella	wet
ω_2	dry, carrying umbrella	dry
$U[\rho(\omega, a)]$	a_1	a_2
ω_1	0	-10
ω_2	0	1

Table: Rewards and utilities.

- $\max_a \min_\omega U = 0$

Uncertain rewards

- Decisions $a \in \mathcal{A}$
- Each choice is called a **reward** $r \in \mathcal{R}$.
- There is a **utility function** $U : \mathcal{R} \rightarrow \mathbb{R}$, assigning values to reward.
- We (weakly) prefer A to B iff $U(A) \geq U(B)$.

Example 7

You are going to work, and it might rain.
What do you do?

- a_1 : Take the umbrella.
- a_2 : Risk it!
- ω_1 : rain
- ω_2 : dry

$\rho(\omega, a)$	a_1	a_2
ω_1	dry, carrying umbrella	wet
ω_2	dry, carrying umbrella	dry
$U[\rho(\omega, a)]$	a_1	a_2
ω_1	0	-10
ω_2	0	1

Table: Rewards and utilities.

- $\max_a \min_\omega U = 0$
- $\min_\omega \max_a U = 0$

Expected utility

$$\mathbb{E}(U | a) = \sum_r U[\rho(\omega, a)] \mathbb{P}(\omega | a)$$

Example 8

You are going to work, and it might rain. The forecast said that the probability of rain (ω_1) was 20%. What do you do?

- a_1 : Take the umbrella.
- a_2 : Risk it!

$\rho(\omega, a)$	a_1	a_2
ω_1	dry, carrying umbrella	wet
ω_2	dry, carrying umbrella	dry
$U[\rho(\omega, a)]$	a_1	a_2
ω_1	0	-10
ω_2	0	1
$\mathbb{E}_P(U a)$	0	-1.2

Table: Rewards, utilities, expected utility for 20% probability of rain.

Bayes decision rules

Consider the case where outcomes are independent of decisions:

$$U(\xi, a) \triangleq \sum_{\mu} U(\mu, a)\xi(\mu)$$

This corresponds e.g. to the case where $\xi(\mu)$ is the belief about an unknown world.

Definition 9 (Bayes utility)

The maximising decision for ξ has an expected utility equal to:

$$U^*(\xi) \triangleq \max_{a \in \mathcal{A}} U(\xi, a). \quad (2.1)$$

The n -meteorologists problem

Exercise 4

- Meteorological models $\mathcal{M} = \{\mu_1, \dots, \mu_n\}$
- Rain predictions at time t : $p_{t,\mu} \triangleq P_\mu(x_t = \text{rain})$.
- Prior probability $\xi(\mu) = 1/n$ for each model.
- Should we take the umbrella?

	M	T	W	T	F	S	S
CNN	0.5	0.6	0.7	0.9	0.5	0.3	0.1
SMHI	0.3	0.7	0.8	0.9	0.5	0.2	0.1
YR	0.6	0.9	0.8	0.5	0.4	0.1	0.1
Rain?	Y	Y	Y	N	Y	N	N

Table: Predictions by three different entities for the probability of rain on a particular day, along with whether or not it actually rained.

The n -meteorologists problem

Exercise 4

	M	T	W	T	F	S	S
CNN	0.5	0.6	0.7	0.9	0.5	0.3	0.1
SMHI	0.3	0.7	0.8	0.9	0.5	0.2	0.1
YR	0.6	0.9	0.8	0.5	0.4	0.1	0.1
Rain?	Y	Y	Y	N	Y	N	N

Table: Predictions by three different entities for the probability of rain on a particular day, along with whether or not it actually rained.

- 1 What is your belief about the quality of each meteorologist after each day?

The n -meteorologists problem

Exercise 4

	M	T	W	T	F	S	S
CNN	0.5	0.6	0.7	0.9	0.5	0.3	0.1
SMHI	0.3	0.7	0.8	0.9	0.5	0.2	0.1
YR	0.6	0.9	0.8	0.5	0.4	0.1	0.1
Rain?	Y	Y	Y	N	Y	N	N

Table: Predictions by three different entities for the probability of rain on a particular day, along with whether or not it actually rained.

- 1 What is your belief about the quality of each meteorologist after each day?
- 2 What is your belief about the probability of rain each day?

$$P_{\xi}(x_t = \text{rain} \mid x_1, x_2, \dots, x_{t-1}) = \sum_{\mu \in \mathcal{M}} P_{\mu}(x_t = \text{rain} \mid x_1, x_2, \dots, x_{t-1}) \xi(\mu \mid x_1, x_2, \dots, x_{t-1})$$

The n -meteorologists problem

Exercise 4

	M	T	W	T	F	S	S
CNN	0.5	0.6	0.7	0.9	0.5	0.3	0.1
SMHI	0.3	0.7	0.8	0.9	0.5	0.2	0.1
YR	0.6	0.9	0.8	0.5	0.4	0.1	0.1
Rain?	Y	Y	Y	N	Y	N	N

Table: Predictions by three different entities for the probability of rain on a particular day, along with whether or not it actually rained.

- 1 What is your belief about the quality of each meteorologist after each day?
- 2 What is your belief about the probability of rain each day?

$$P_{\xi}(x_t = \text{rain} \mid x_1, x_2, \dots, x_{t-1}) = \sum_{\mu \in \mathcal{M}} P_{\mu}(x_t = \text{rain} \mid x_1, x_2, \dots, x_{t-1}) \xi(\mu \mid x_1, x_2, \dots, x_{t-1})$$

- 3 Assume you can decide whether or not to go running each day. If you go running and it does not rain, your utility is 1. If it rains, it's -10. If you don't go running, your utility is 0. What is the decision maximising utility in expectation (with respect to the posterior) each day?

Deciding a class given a model

- Features $x_t \in \mathcal{X}$.
- Label $y_t \in \mathcal{Y}$.
- Decisions $a_t \in \mathcal{A}$.
- Decision rule $\pi(a_t | x_t)$ assigns probabilities to actions.

Standard classification problem

$$\mathcal{A} = \mathcal{Y}, \quad U(a, y) = \mathbb{I}\{a = y\}$$

Exercise 5

If we have a model $P_\mu(y_t | x_t)$, and a suitable U , what is the optimal decision to make?

Deciding a class given a model

- Features $x_t \in \mathcal{X}$.
- Label $y_t \in \mathcal{Y}$.
- Decisions $a_t \in \mathcal{A}$.
- Decision rule $\pi(a_t | x_t)$ assigns probabilities to actions.

Standard classification problem

$$\mathcal{A} = \mathcal{Y}, \quad U(a, y) = \mathbb{I}\{a = y\}$$

Exercise 5

If we have a model $P_\mu(y_t | x_t)$, and a suitable U , what is the optimal decision to make?

$$a_t \in \arg \max_{a \in \mathcal{A}} \sum_y P_\mu(y_t = y | x_t) U(a, y)$$

Deciding a class given a model

- Features $x_t \in \mathcal{X}$.
- Label $y_t \in \mathcal{Y}$.
- Decisions $a_t \in \mathcal{A}$.
- Decision rule $\pi(a_t | x_t)$ assigns probabilities to actions.

Standard classification problem

$$\mathcal{A} = \mathcal{Y}, \quad U(a, y) = \mathbb{I}\{a = y\}$$

Exercise 5

If we have a model $P_\mu(y_t | x_t)$, and a suitable U , what is the optimal decision to make?

$$a_t \in \arg \max_{a \in \mathcal{A}} \sum_y P_\mu(y_t = y | x_t) U(a, y)$$

For standard classification,

$$a_t \in \arg \max_{a \in \mathcal{A}} P_\mu(y_t = a | x_t)$$

Deciding the class given a model family

- Training data $D_T = \{(x_i, y_i) \mid i = 1, \dots, T\}$
- Models $\{P_\mu \mid \mu \in \mathcal{M}\}$.
- Prior ξ on \mathcal{M} .

Posterior over classification models

$$\xi(\mu \mid D_T) = \frac{P_\mu(y_1, \dots, y_T \mid x_1, \dots, x_T)\xi(\mu)}{\sum_{\mu' \in \mathcal{M}} P_{\mu'}(y_1, \dots, y_T \mid x_1, \dots, x_T)\xi(\mu')}$$

Deciding the class given a model family

- Training data $D_T = \{(x_i, y_i) \mid i = 1, \dots, T\}$
- Models $\{P_\mu \mid \mu \in \mathcal{M}\}$.
- Prior ξ on \mathcal{M} .

Posterior over classification models

$$\xi(\mu \mid D_T) = \frac{P_\mu(y_1, \dots, y_T \mid x_1, \dots, x_T) \xi(\mu)}{\sum_{\mu' \in \mathcal{M}} P_{\mu'}(y_1, \dots, y_T \mid x_1, \dots, x_T) \xi(\mu')}$$

If not dealing with time-series data, we assume independence between x_t :

$$P_\mu(y_1, \dots, y_T \mid x_1, \dots, x_T) = \prod_{i=1}^T P_\mu(y_i \mid x_i)$$

Deciding the class given a model family

- Training data $D_T = \{(x_i, y_i) \mid i = 1, \dots, T\}$
- Models $\{P_\mu \mid \mu \in \mathcal{M}\}$.
- Prior ξ on \mathcal{M} .

Posterior over classification models

$$\xi(\mu \mid D_T) = \frac{P_\mu(y_1, \dots, y_T \mid x_1, \dots, x_T)\xi(\mu)}{\sum_{\mu' \in \mathcal{M}} P_{\mu'}(y_1, \dots, y_T \mid x_1, \dots, x_T)\xi(\mu')}$$

The Bayes rule for maximising $\mathbb{E}_\xi(U \mid a, x_t, D_T)$

The decision rule simply chooses the action:

$$a_t \in \arg \max_{a \in \mathcal{A}} \sum_y \sum_{\mu \in \mathcal{M}} P_\mu(y_t = y \mid x_t)\xi(\mu \mid D_T)U(a, y) \quad (3.1)$$

Deciding the class given a model family

- Training data $D_T = \{(x_i, y_i) \mid i = 1, \dots, T\}$
- Models $\{P_\mu \mid \mu \in \mathcal{M}\}$.
- Prior ξ on \mathcal{M} .

Posterior over classification models

$$\xi(\mu \mid D_T) = \frac{P_\mu(y_1, \dots, y_T \mid x_1, \dots, x_T) \xi(\mu)}{\sum_{\mu' \in \mathcal{M}} P_{\mu'}(y_1, \dots, y_T \mid x_1, \dots, x_T) \xi(\mu')}$$

The Bayes rule for maximising $\mathbb{E}_\xi(U \mid a, x_t, D_T)$

The decision rule simply chooses the action:

$$a_t \in \arg \max_{a \in \mathcal{A}} \sum_y \sum_{\mu \in \mathcal{M}} P_\mu(y_t = y \mid x_t) \xi(\mu \mid D_T) U(a, y) \quad (3.1)$$

We can rewrite this by calculating the posterior marginal label probability

$$\mathbb{P}_{\xi \mid D_T}(y_t \mid x_t) \triangleq \mathbb{P}_\xi(y_t \mid x_t, D_T) = \sum_{\mu \in \mathcal{M}} P_\mu(y_t \mid x_t) \xi(\mu \mid D_T).$$

Deciding the class given a model family

- Training data $D_T = \{(x_i, y_i) \mid i = 1, \dots, T\}$
- Models $\{P_\mu \mid \mu \in \mathcal{M}\}$.
- Prior ξ on \mathcal{M} .

Posterior over classification models

$$\xi(\mu \mid D_T) = \frac{P_\mu(y_1, \dots, y_T \mid x_1, \dots, x_T) \xi(\mu)}{\sum_{\mu' \in \mathcal{M}} P_{\mu'}(y_1, \dots, y_T \mid x_1, \dots, x_T) \xi(\mu')}$$

The Bayes rule for maximising $\mathbb{E}_\xi(U \mid a, x_t, D_T)$

The decision rule simply chooses the action:

$$a_t \in \arg \max_{a \in \mathcal{A}} \sum_y \sum_{\mu \in \mathcal{M}} P_\mu(y_t = y \mid x_t) \xi(\mu \mid D_T) U(a, y) \quad (3.1)$$

$$= \arg \max_{a \in \mathcal{A}} \sum_y \mathbb{P}_{\xi \mid D_T}(y_t \mid x_t) U(a, y) \quad (3.2)$$

We can rewrite this by calculating the posterior marginal label probability

$$\mathbb{P}_{\xi \mid D_T}(y_t \mid x_t) \triangleq \mathbb{P}_\xi(y_t \mid x_t, D_T) = \sum_{\mu \in \mathcal{M}} P_\mu(y_t \mid x_t) \xi(\mu \mid D_T).$$

Approximating the model

Full Bayesian approach for infinite \mathcal{M}

Here ξ can be a probability density function and

$$\xi(\mu | D_T) = P_\mu(D_T)\xi(\mu) / \mathbb{P}_\xi(D_T), \quad \mathbb{P}_\xi(D_T) = \int_{\mathcal{M}} P_\mu(D_T)\xi(\mu) d,$$

can be hard to calculate.

Approximating the model

Full Bayesian approach for infinite \mathcal{M}

Here ξ can be a probability density function and

$$\xi(\mu \mid D_T) = P_\mu(D_T)\xi(\mu) / \mathbb{P}_\xi(D_T), \quad \mathbb{P}_\xi(D_T) = \int_{\mathcal{M}} P_\mu(D_T)\xi(\mu) d,$$

can be hard to calculate.

Maximum a posteriori model

We only choose a single model through the following optimisation:

$$\mu_{\text{MAP}}(\xi, D_T) = \arg \max_{\mu \in \mathcal{M}} P_\mu(D_T)\xi(\mu)$$

Approximating the model

Full Bayesian approach for infinite \mathcal{M}

Here ξ can be a probability density function and

$$\xi(\mu | D_T) = P_\mu(D_T)\xi(\mu) / \mathbb{P}_\xi(D_T), \quad \mathbb{P}_\xi(D_T) = \int_{\mathcal{M}} P_\mu(D_T)\xi(\mu) d,$$

can be hard to calculate.

Maximum a posteriori model

We only choose a single model through the following optimisation:

$$\mu_{\text{MAP}}(\xi, D_T) = \arg \max_{\mu \in \mathcal{M}} \overbrace{\ln P_\mu(D_T)}^{\text{goodness of fit}} + \underbrace{\ln \xi(\mu)}_{\text{regulariser}} .$$

Learning outcomes

Understanding

- Preferences, utilities and the expected utility principle.
- Hypothesis testing and classification as decision problems.
- How to interpret p -values Bayesian tests.
- The MAP approximation to full Bayesian inference.

Skills

- Being able to implement an optimal decision rule for a given utility and probability.
- Being able to construct a simple null hypothesis test.

Reflection

- When would expected utility maximisation not be a good idea?
- What does a p value represent when you see it in a paper?
- Can we prevent high false discovery rates when using p values?
- When is the MAP approximation good?

Simple hypothesis testing

The simple hypothesis test as a decision problem

- $\mathcal{M} = \{\mu_0, \mu_1\}$
- a_0 : Accept model μ_0
- a_1 : Accept model μ_1

U	μ_0	μ_1
a_0	1	0
a_1	0	1

Table: Example utility function for simple hypothesis tests.

Example 10 (Continuation of the medium example)

- μ_1 : that John is a medium.
- μ_0 : that John is not a medium.

$$\mathbb{E}_\xi(U \mid a_0) = 1 \times \xi(\mu_0 \mid \mathbf{x}) + 0 \times \xi(\mu_1 \mid \mathbf{x}), \quad \mathbb{E}_\xi(U \mid a_1) = 0 \times \xi(\mu_0 \mid \mathbf{x}) + 1 \times \xi(\mu_1 \mid \mathbf{x})$$

Null hypothesis test

Many times, there is only one model under consideration, μ_0 , the so-called **null hypothesis**.

The null hypothesis test as a decision problem

- a_0 : Accept model μ_0
- a_1 : Reject model μ_0

Example 11

Construction of the test for the medium

Null hypothesis test

Many times, there is only one model under consideration, μ_0 , the so-called **null hypothesis**.

The null hypothesis test as a decision problem

- a_0 : Accept model μ_0
- a_1 : Reject model μ_0

Example 11

Construction of the test for the medium

- μ_0 is simply the *Bernoulli*(1/2) model: responses are by chance.

Null hypothesis test

Many times, there is only one model under consideration, μ_0 , the so-called **null hypothesis**.

The null hypothesis test as a decision problem

- a_0 : Accept model μ_0
- a_1 : Reject model μ_0

Example 11

Construction of the test for the medium

- μ_0 is simply the *Bernoulli*(1/2) model: responses are by chance.
- We need to design a policy $\pi(a | \mathbf{x})$ that accepts or rejects depending on the data.

Null hypothesis test

Many times, there is only one model under consideration, μ_0 , the so-called **null hypothesis**.

The null hypothesis test as a decision problem

- a_0 : Accept model μ_0
- a_1 : Reject model μ_0

Example 11

Construction of the test for the medium

- μ_0 is simply the *Bernoulli*(1/2) model: responses are by chance.
- We need to design a policy $\pi(a | \mathbf{x})$ that accepts or rejects depending on the data.
- Since there is no alternative model, we can only construct this policy according to its properties when μ_0 is true.

Null hypothesis test

Many times, there is only one model under consideration, μ_0 , the so-called **null hypothesis**.

The null hypothesis test as a decision problem

- a_0 : Accept model μ_0
- a_1 : Reject model μ_0

Example 11

Construction of the test for the medium

- μ_0 is simply the *Bernoulli*(1/2) model: responses are by chance.
- We need to design a policy $\pi(a | \mathbf{x})$ that accepts or rejects depending on the data.
- Since there is no alternative model, we can only construct this policy according to its properties when μ_0 is true.
- In particular, we can fix a policy that only chooses a_1 when μ_0 is true a proportion δ of the time.

Null hypothesis test

Many times, there is only one model under consideration, μ_0 , the so-called **null hypothesis**.

The null hypothesis test as a decision problem

- a_0 : Accept model μ_0
- a_1 : Reject model μ_0

Example 11

Construction of the test for the medium

- μ_0 is simply the *Bernoulli*(1/2) model: responses are by chance.
- We need to design a policy $\pi(a | \mathbf{x})$ that accepts or rejects depending on the data.
- Since there is no alternative model, we can only construct this policy according to its properties when μ_0 is true.
- In particular, we can fix a policy that only chooses a_1 when μ_0 is true a proportion δ of the time.
- This can be done by constructing a threshold test from the inverse-CDF.

Using p -values to construct statistical tests

Definition 12 (Null statistical test)

The statistic $f : \mathcal{X} \rightarrow [0, 1]$ is designed to have the property:

$$P_{\mu_0}(\{x \mid f(x) \leq \delta\}) = \delta.$$

If our decision rule is:

$$\pi(a \mid x) = \begin{cases} a_0, & f(x) \leq \delta \\ a_1, & f(x) > \delta, \end{cases}$$

the probability of rejecting the null hypothesis when it is true is exactly δ .

The value of the statistic $f(x)$, otherwise known as the p -value, is uninformative.

Issues with p -values

- They only measure quality of fit **on the data**.
- Not robust to model misspecification.
- They ignore effect sizes.
- They do not consider prior information.
- They do not represent the probability of having made an error.
- The null-rejection error probability is the same irrespective of the amount of data (by design).

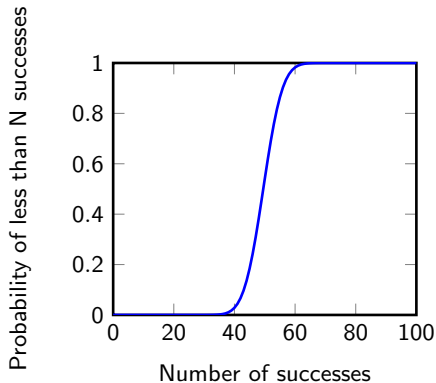
p -values for the medium example

p -values for the medium example

- μ_0 is simply the *Bernoulli*(1/2) model: responses are by chance.

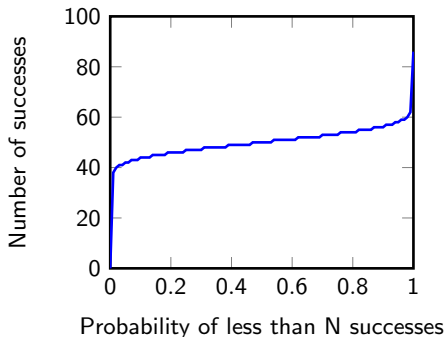
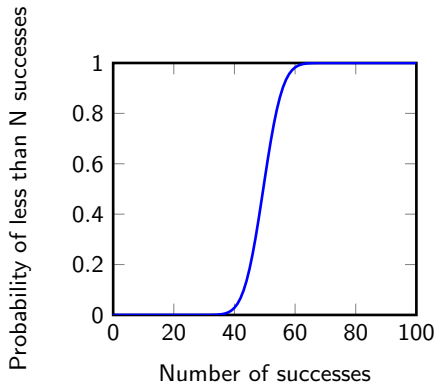
p -values for the medium example

- μ_0 is simply the *Bernoulli*(1/2) model: responses are by chance.
- CDF: $P_{\mu_0}(N \leq n \mid K = 100)$



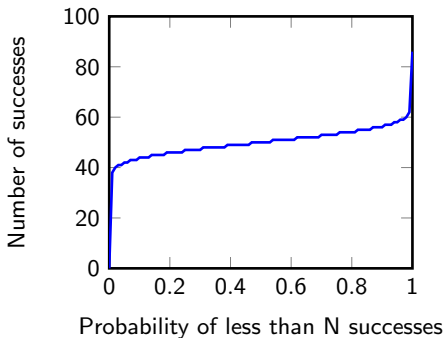
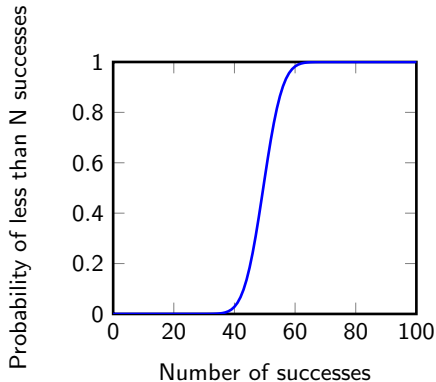
p -values for the medium example

- μ_0 is simply the *Bernoulli*(1/2) model: responses are by chance.
- CDF: $P_{\mu_0}(N \leq n \mid K = 100)$
- ICDF: the number of successes that will happen with probability at least δ



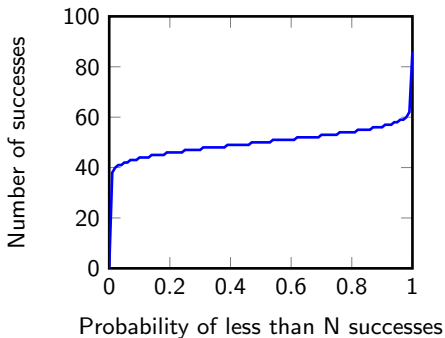
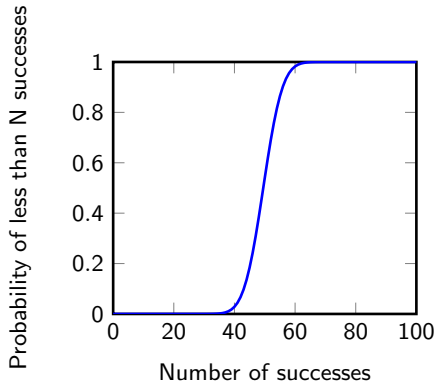
p -values for the medium example

- μ_0 is simply the *Bernoulli*(1/2) model: responses are by chance.
- CDF: $P_{\mu_0}(N \leq n \mid K = 100)$
- ICDF: the number of successes that will happen with probability at least δ
- e.g. we'll get at most 50 successes a proportion $\delta = 1/2$ of the time.



p -values for the medium example

- μ_0 is simply the *Bernoulli*(1/2) model: responses are by chance.
- CDF: $P_{\mu_0}(N \leq n \mid K = 100)$
- ICDF: the number of successes that will happen with probability at least δ
- e.g. we'll get at most 50 successes a proportion $\delta = 1/2$ of the time.
- Using the (inverse) CDF we can construct a policy π that selects a_1 when μ_0 is true only a δ portion of the time, for any choice of δ .



Building a test

The test statistic

We want the test to reflect that we don't have a significant number of failures.

$$f(x) = 1 - \text{binocdf}\left(\sum_{t=1}^n x_t, n, 0.5\right)$$

What $f(x)$ is and is not

- It is a **statistic** which is $\leq \delta$ a δ portion of the time when μ_0 is true.
- It is **not** the probability of observing x under μ_0 .
- It is **not** the probability of μ_0 given x .

Exercise 6

- Let us throw a coin 8 times, and try and predict the outcome.

Exercise 6

- Let us throw a coin 8 times, and try and predict the outcome.
- Select a p -value threshold so that $\delta = 0.05$. For 8 throws, this corresponds to

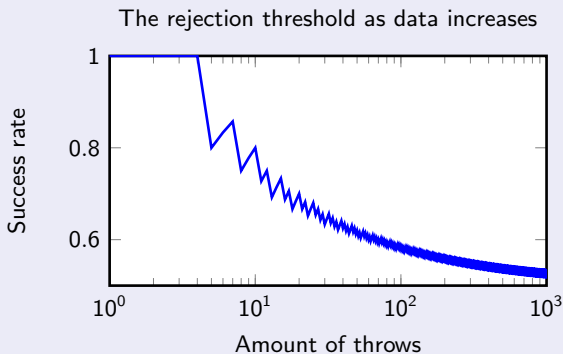


Figure: Here we see how the rejection threshold, in terms of the success rate, changes with the number of throws to achieve an error rate of $\delta = 0.05$.

Exercise 6

- Let us throw a coin 8 times, and try and predict the outcome.
- Select a p -value threshold so that $\delta = 0.05$. For 8 throws, this corresponds to > 6 successes or $\geq 87.5\%$ success rate.
- Let's calculate the p -value for each one of you

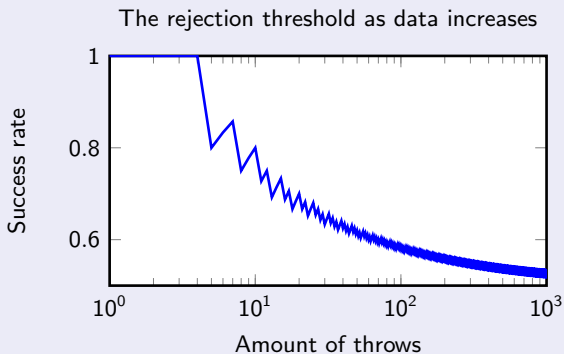


Figure: Here we see how the rejection threshold, in terms of the success rate, changes with the number of throws to achieve an error rate of $\delta = 0.05$.

Exercise 6

- Let us throw a coin 8 times, and try and predict the outcome.
- Select a p -value threshold so that $\delta = 0.05$. For 8 throws, this corresponds to > 6 successes or $\geq 87.5\%$ success rate.
- Let's calculate the p -value for each one of you
- What is the rejection performance of the test?

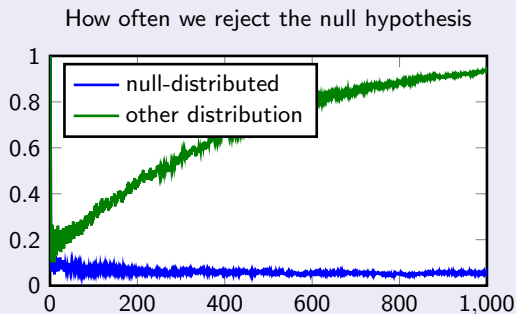


Figure: Here we see the rejection rate of the null hypothesis (μ_0) for two cases. Firstly, for the case when μ_0 is true. Secondly, when the data is generated from $Bernoulli(0.55)$.

Statistical power and false discovery.

Beyond not rejecting the null when it's true, we also want:

- High power: Rejecting the null when it is false.
- Low false discovery rate: Accepting the null when it is true.

Power

The power depends on what hypothesis we use as an alternative.

False discovery rate

False discovery depends on how likely it is **a priori** that the null is false.

The Bayesian version of the test

Example 13

- 1 Set $U(a_i, \mu_j) = \mathbb{I}\{i = j\}$.
- 2 Set $\xi(\mu_i) = 1/2$.
- 3 μ_0 : *Bernoulli*(1/2).
- 4 μ_1 : *Bernoulli*(θ), $\theta \sim \text{Unif}([0, 1])$.
- 5 Calculate $\xi(\mu | x)$.
- 6 Choose a_i , where $i = \arg \max_j \xi(\mu_j | x)$.

Bayesian model averaging for the alternative model μ_1

$$P_{\mu_1}(x) = \int_{\Theta} B_{\theta}(x) d\beta(\theta) \quad (3.3)$$

$$\xi(\mu_0 | x) = \frac{P_{\mu_0}(x)\xi(\mu_0)}{P_{\mu_0}(x)\xi(\mu_0) + P_{\mu_1}(x)\xi(\mu_1)} \quad (3.4)$$

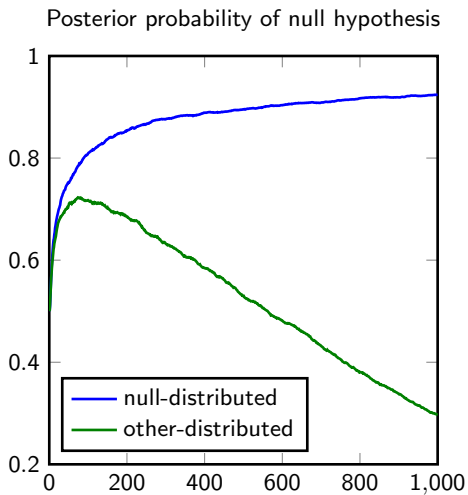


Figure: Here we see the convergence of the posterior probability.

Rejection of null hypothesis for Bernoulli(0.5)

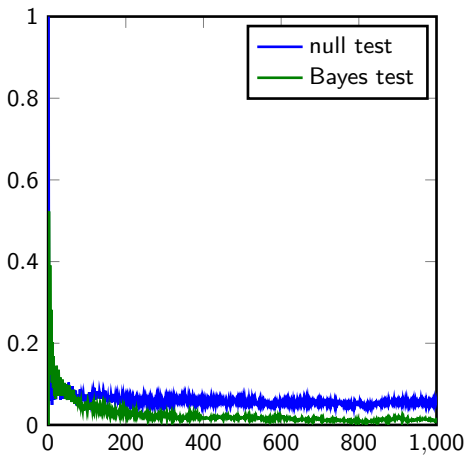


Figure: Comparison of the rejection probability for the null and the Bayesian test when μ_0 is true.

Rejection of null hypothesis for Bernoulli(0.55)

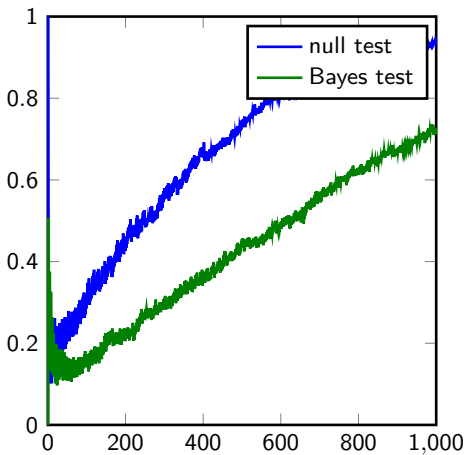


Figure: Comparison of the rejection probability for the null and the Bayesian test when μ_1 is true.

Further reading

Points of significance (Nature Methods)

- Importance of being uncertain <https://www.nature.com/articles/nmeth.2613>
- Error bars <https://www.nature.com/articles/nmeth.2659>
- P values and the search for significance
<https://www.nature.com/articles/nmeth.4120>
- Bayes' theorem <https://www.nature.com/articles/nmeth.3335>
- Sampling distributions and the bootstrap
<https://www.nature.com/articles/nmeth.3414>

- 1 Beliefs and probabilities
- 2 Hierarchies of decision making problems
- 3 Formalising Classification problems
- 4 Classification with stochastic gradient descent**
 - Neural network models

Classification as an optimisation problem.

The μ -optimal classifier

$$\max_{\theta \in \Theta} f(\pi_{\theta}, \mu, U),$$

$$f(\pi_{\theta}, \mu, U) \triangleq \mathbb{E}_{\mu}^{\pi_{\theta}}(U) \quad (4.1)$$

$$f(\pi_{\theta}, \mu, U) = \sum_{x, y, a} U(a, y) \pi_{\theta}(a | x) P_{\mu}(y | x) P_{\mu}(x) \quad (4.2)$$

$$\approx \sum_{t=1}^T \sum_{a_t} U(a_t, y_t) \pi_{\theta}(a_t | x_t), \quad (x_t, y_t)_{t=1}^T \sim P_{\mu}. \quad (4.3)$$

Bayesian inference for Bernoulli distributions

Estimating a coin's bias

A fair coin comes heads 50% of the time. We want to test an unknown coin, which we think may not be completely fair.

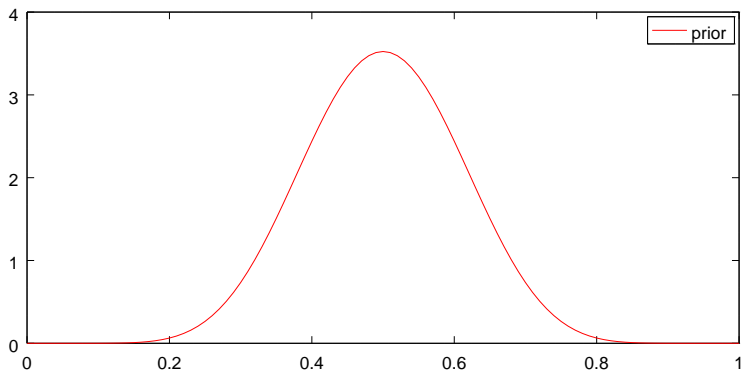


Figure: Prior belief ξ about the coin bias θ .

Bayesian inference for Bernoulli distributions

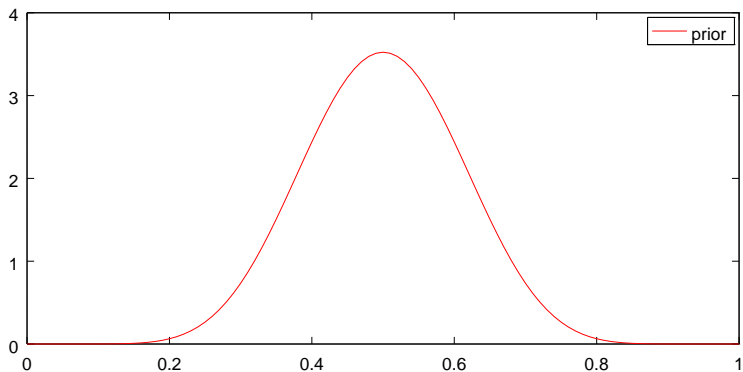


Figure: Prior belief ξ about the coin bias θ .

For a sequence of throws $x_t \in \{0, 1\}$,

$$P_{\theta}(x) \propto \prod_t \theta^{x_t} (1 - \theta)^{1 - x_t} = \theta^{\#\text{Heads}} (1 - \theta)^{\#\text{Tails}}$$

Bayesian inference for Bernoulli distributions

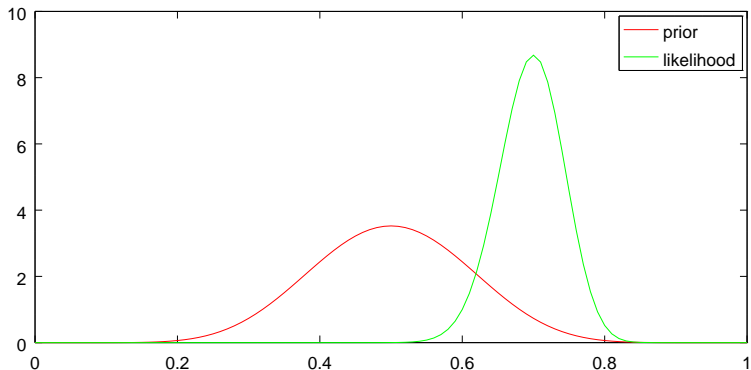


Figure: Prior belief ξ about the coin bias θ and likelihood of θ for the data.

Say we throw the coin 100 times and obtain 70 heads. Then we plot the **likelihood** $P_{\theta}(x)$ of different models.

Bayesian inference for Bernoulli distributions

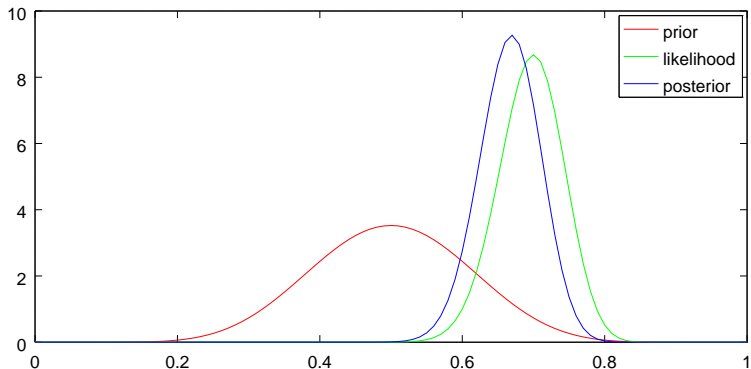


Figure: Prior belief $\xi(\theta)$ about the coin bias θ , likelihood of θ for the data, and posterior belief $\xi(\theta | x)$

From these, we calculate a **posterior** distribution over the correct models. This represents our conclusion given our prior and the data.

Stochastic gradient methods

Gradient ascent

$$\theta_{i+1} = \theta_i + \alpha \nabla_{\theta} g(\theta_i).$$

Stochastic gradient ascent

$$g(\theta) = \int_{\mathcal{M}} f(\theta, \mu) d\xi(\mu)$$
$$\theta_{i+1} = \theta_i + \alpha \nabla_{\theta} f(\theta_i, \mu_i), \quad \mu_i \sim \xi.$$

Two views of neural networks

Neural network classification model $P_{\theta}(\mathbf{y} \mid \mathbf{x}_t)$



Objective: Find the best model for D_T .

Neural network classification policy $\pi(a_t \mid \mathbf{x}_t)$



Objective: Find the best policy for $U(a, \mathbf{x})$.

Two views of neural networks

Neural network classification model $P_{\theta}(\mathbf{y} \mid \mathbf{x}_t)$ 

Objective: Find the best model for D_T .

Neural network classification policy $\pi(a_t \mid \mathbf{x}_t)$ 

Objective: Find the best policy for $U(a, \mathbf{x})$.

Difference between the two views

- We can use standard probabilistic methods for P .
- Finding the optimal π is an optimisation problem.

Linear networks and the perceptron algorithm



Figure: Abstract graphical model for a neural network

Definition 14 (Linear classifier)

$$\Theta = [\theta_1 \quad \cdots \quad \theta_C] = \begin{bmatrix} \theta_{1,1} & \cdots & \theta_{1,C} \\ \vdots & \ddots & \vdots \\ \theta_N & \cdots & \theta_{N,C} \end{bmatrix}$$

$$\pi_{\Theta}(a \mid \mathbf{x}) = \exp(\theta_a^{\top} \mathbf{x}) / \sum_{a'} \exp(\theta_{a'}^{\top} \mathbf{x})$$

Linear networks and the perceptron algorithm

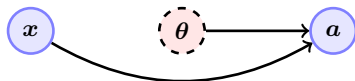


Figure: Abstract graphical model for a neural network

Definition 14 (Linear classifier)

$$\Theta = [\theta_1 \quad \cdots \quad \theta_C] = \begin{bmatrix} \theta_{1,1} & \cdots & \theta_{1,C} \\ \vdots & \ddots & \vdots \\ \theta_N & \cdots & \theta_{N,C} \end{bmatrix}$$

$$\pi_{\Theta}(a \mid \mathbf{x}) = \exp(\theta_a^{\top} \mathbf{x}) / \sum_{a'} \exp(\theta_{a'}^{\top} \mathbf{x})$$

Linear networks and the perceptron algorithm

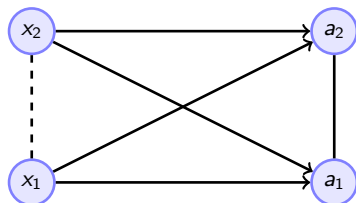


Figure: Graphical model for a linear neural network.

Definition 14 (Linear classifier)

$$\Theta = [\theta_1 \quad \cdots \quad \theta_C] = \begin{bmatrix} \theta_{1,1} & \cdots & \theta_{1,C} \\ \vdots & \ddots & \vdots \\ \theta_N & \cdots & \theta_{N,C} \end{bmatrix}$$

$$\pi_{\Theta}(a \mid \mathbf{x}) = \exp(\theta_a^{\top} \mathbf{x}) / \sum_{a'} \exp(\theta_{a'}^{\top} \mathbf{x})$$

Linear networks and the perceptron algorithm

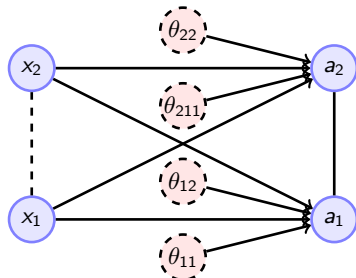


Figure: Graphical model for a linear neural network.

Definition 14 (Linear classifier)

$$\Theta = [\theta_1 \quad \cdots \quad \theta_C] = \begin{bmatrix} \theta_{1,1} & \cdots & \theta_{1,C} \\ \vdots & \ddots & \vdots \\ \theta_N & \cdots & \theta_{N,C} \end{bmatrix}$$

$$\pi_{\Theta}(a \mid \mathbf{x}) = \exp(\theta_a^{\top} \mathbf{x}) / \sum_{a'} \exp(\theta_{a'}^{\top} \mathbf{x})$$

Linear networks and the perceptron algorithm

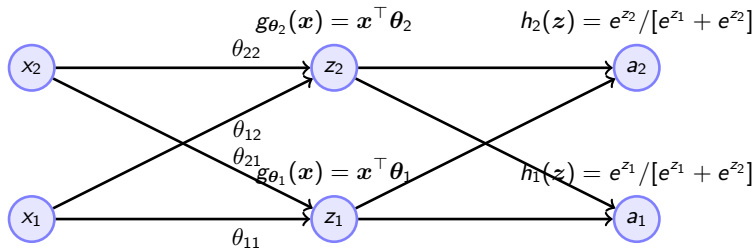


Figure: Architectural view of a linear neural network.

Definition 14 (Linear classifier)

$$\Theta = [\boldsymbol{\theta}_1 \quad \cdots \quad \boldsymbol{\theta}_C] = \begin{bmatrix} \theta_{1,1} & \cdots & \theta_{1,C} \\ \vdots & \ddots & \vdots \\ \theta_{N,1} & \cdots & \theta_{N,C} \end{bmatrix}$$

$$\pi_{\Theta}(a | \mathbf{x}) = \exp(\boldsymbol{\theta}_a^\top \mathbf{x}) / \sum_{a'} \exp(\boldsymbol{\theta}_{a'}^\top \mathbf{x})$$

Gradient ascent for a matrix U

$$\max_{\theta} \sum_{t=1}^T \sum_{a_t} U(a_t, y_t) \pi_{\theta}(a_t | x_t) \quad (\text{objective})$$

$$\nabla_{\theta} \sum_{t=1}^T \sum_{a_t} U(a_t, y_t) \pi_{\theta}(a_t | x_t) \quad (\text{gradient})$$

$$= \sum_{t=1}^T \sum_{a_t} U(a_t, y_t) \nabla_{\theta} \pi_{\theta}(a_t | x_t) \quad (4.4)$$

Chain Rule of Differentiation

$$f(z), z = g(x),$$

$$\frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx} \quad (\text{scalar version})$$

$$\nabla_{\theta} \pi = \nabla_g \pi \nabla_{\theta} g \quad (\text{vector version})$$

Learning outcomes

Understanding

- Classification as an optimisation problem.
- (Stochastic) gradient methods and the chain rule.
- Neural networks as probability models or classification policies.
- Linear neural networks.
- Nonlinear network architectures.

Skills

- Using a standard NN class in python.

Reflection

- How useful is the ability to have multiple non-linear layers in a neural network.
- How rich is the representational power of neural networks?
- Is there anything special about neural networks other than their allusions to biology?