Probabilistic & Unsupervised Learning

Introduction and Foundations

Maneesh Sahani maneesh@gatsby.ucl.ac.uk

Gatsby Computational Neuroscience Unit, and MSc ML/CSML, Dept Computer Science University College London

Term 1, Autumn 2016



Jan Steen

Not just remembering:



Jan Steen

Not just remembering:

Systematising (noisy) observations: discovering structure.



Jan Steen

Not just remembering:

- Systematising (noisy) observations: discovering structure.
- Predicting new outcomes: generalising.



Jan Steen

Not just remembering:

- Systematising (noisy) observations: discovering structure.
- Predicting new outcomes: generalising.
- Choosing actions wisely.

Systematising (noisy) observations: discovering structure.

Predicting new outcomes: generalising.

Choosing actions wisely.

Systematising (noisy) observations: discovering structure.

• Unsupervised learning. Observe (sensory) input alone:

 $x_1, x_2, x_3, x_4, \ldots$

Describe pattern of data [p(x)], identify and extract underlying structural variables [$x_i \rightarrow y_i$].

Predicting new outcomes: generalising.

Choosing actions wisely.

Systematising (noisy) observations: discovering structure.

Unsupervised learning. Observe (sensory) input alone:

 $x_1, x_2, x_3, x_4, \ldots$

Describe pattern of data [p(x)], identify and extract underlying structural variables [$x_i \rightarrow y_i$].

- Predicting new outcomes: generalising.
 - Supervised learning. Observe input/output pairs ("teaching"):

 $(x_1, y_1), (x_2, y_2), (x_3, y_3), (x_4, y_4), \ldots$

Predict the correct y^* for new test input x^* .

Choosing actions wisely.

Systematising (noisy) observations: discovering structure.

Unsupervised learning. Observe (sensory) input alone:

 $x_1, x_2, x_3, x_4, \ldots$

Describe pattern of data [p(x)], identify and extract underlying structural variables $[x_i \rightarrow y_i]$.

- Predicting new outcomes: generalising.
 - Supervised learning. Observe input/output pairs ("teaching"):

 $(x_1, y_1), (x_2, y_2), (x_3, y_3), (x_4, y_4), \ldots$

Predict the correct y^* for new test input x^* .

- Choosing actions wisely.
 - Reinforcement learning. Rewards or payoffs (and possibly also inputs) depend on actions:

 $x_1 : a_1 \rightarrow r_1, x_2 : a_2 \rightarrow r_2, x_3 : a_3 \rightarrow r_3 \dots$

Find a policy for action choice that maximises payoff.

Unsupervised Learning

Find underlying structure:

- separate generating processes (clusters)
- reduced dimensionality representations
- good explanations (causes) of the data
- modelling the data density



- structure discovery, science
- data compression
- outlier detection
- input to supervised/reinforcement algorithms (causes may be more simply related to outputs or rewards)
- a theory of biological learning and perception



Supervised learning

Two main examples:



Discrete (class label) outputs.

But also: ranks, relationships, trees etc.

Variants may relate to unsupervised learning:

- semi-supervised learning (most *x* unlabelled; assumes structure of $\{x\}$ and relationship $x \rightarrow y$ are linked).
- multitask (transfer) learning (predict different y in different contexts; assumes links between structure of relationships).



Continuous-values outputs.

Data are generated by random and/or unknown processes.

Data are generated by random and/or unknown processes.

Our approach to learning starts with a probabilistic model of data production:

 $P(\text{data}|\text{parameters}) = P(x|\theta) \text{ or } P(y|x,\theta)$

Data are generated by random and/or unknown processes.

Our approach to learning starts with a probabilistic model of data production:

 $P(\text{data}|\text{parameters}) = P(x|\theta) \text{ or } P(y|x,\theta)$

This is the generative model or likelihood.

The probabilistic model can be used to

Data are generated by random and/or unknown processes.

Our approach to learning starts with a probabilistic model of data production:

 $P(\text{data}|\text{parameters}) = P(x|\theta) \text{ or } P(y|x,\theta)$

- The probabilistic model can be used to
 - make inferences about missing inputs

Data are generated by random and/or unknown processes.

Our approach to learning starts with a probabilistic model of data production:

 $P(\text{data}|\text{parameters}) = P(x|\theta) \text{ or } P(y|x,\theta)$

- The probabilistic model can be used to
 - make inferences about missing inputs
 - generate predictions/fantasies/imagery

Data are generated by random and/or unknown processes.

Our approach to learning starts with a probabilistic model of data production:

 $P(\text{data}|\text{parameters}) = P(x|\theta) \text{ or } P(y|x,\theta)$

- The probabilistic model can be used to
 - make inferences about missing inputs
 - generate predictions/fantasies/imagery
 - make predictions or decisions which minimise expected loss

Data are generated by random and/or unknown processes.

Our approach to learning starts with a probabilistic model of data production:

 $P(\text{data}|\text{parameters}) = P(x|\theta) \text{ or } P(y|x,\theta)$

- The probabilistic model can be used to
 - make inferences about missing inputs
 - generate predictions/fantasies/imagery
 - make predictions or decisions which minimise expected loss
 - communicate the data in an efficient way

Data are generated by random and/or unknown processes.

Our approach to learning starts with a probabilistic model of data production:

 $P(\text{data}|\text{parameters}) = P(x|\theta) \text{ or } P(y|x,\theta)$

- The probabilistic model can be used to
 - make inferences about missing inputs
 - generate predictions/fantasies/imagery
 - make predictions or decisions which minimise expected loss
 - communicate the data in an efficient way
- Probabilistic modelling is often equivalent to other views of learning:

Data are generated by random and/or unknown processes.

Our approach to learning starts with a probabilistic model of data production:

 $P(\text{data}|\text{parameters}) = P(x|\theta) \text{ or } P(y|x,\theta)$

- The probabilistic model can be used to
 - make inferences about missing inputs
 - generate predictions/fantasies/imagery
 - make predictions or decisions which minimise expected loss
 - communicate the data in an efficient way
- Probabilistic modelling is often equivalent to other views of learning:
 - information theoretic: finding compact representations of the data

Data are generated by random and/or unknown processes.

Our approach to learning starts with a probabilistic model of data production:

 $P(\text{data}|\text{parameters}) = P(x|\theta) \text{ or } P(y|x,\theta)$

- The probabilistic model can be used to
 - make inferences about missing inputs
 - generate predictions/fantasies/imagery
 - make predictions or decisions which minimise expected loss
 - communicate the data in an efficient way
- Probabilistic modelling is often equivalent to other views of learning:
 - information theoretic: finding compact representations of the data
 - physical analogies: minimising (free) energy of a corresponding statistical mechanical system

Data are generated by random and/or unknown processes.

Our approach to learning starts with a probabilistic model of data production:

 $P(\text{data}|\text{parameters}) = P(x|\theta) \text{ or } P(y|x,\theta)$

- The probabilistic model can be used to
 - make inferences about missing inputs
 - generate predictions/fantasies/imagery
 - make predictions or decisions which minimise expected loss
 - communicate the data in an efficient way
- Probabilistic modelling is often equivalent to other views of learning:
 - information theoretic: finding compact representations of the data
 - physical analogies: minimising (free) energy of a corresponding statistical mechanical system
 - structural risk: compensate for overconfidence in powerful models

Data are generated by random and/or unknown processes.

Our approach to learning starts with a probabilistic model of data production:

 $P(\text{data}|\text{parameters}) = P(x|\theta) \text{ or } P(y|x,\theta)$

This is the generative model or likelihood.

- The probabilistic model can be used to
 - make inferences about missing inputs
 - generate predictions/fantasies/imagery
 - make predictions or decisions which minimise expected loss
 - communicate the data in an efficient way
- Probabilistic modelling is often equivalent to other views of learning:
 - information theoretic: finding compact representations of the data
 - physical analogies: minimising (free) energy of a corresponding statistical mechanical system
 - structural risk: compensate for overconfidence in powerful models

The calculus of probabilities naturally handles randomness. It is also the right way to reason about unknown values.

Representing beliefs

Let b(x) represent our strength of belief in (plausibility of) proposition x:

$0 \leq b(x) \leq 1$		
b(x)=0	х	is definitely not true
b(x) = 1	х	is definitely true
b(x y)	str	ength of belief that x is true given that we know y is true

Cox Axioms (Desiderata):

- ▶ Let b(x) be real. As b(x) increases, $b(\neg x)$ decreases, and so the function mapping $b(x) \leftrightarrow b(\neg x)$ is monotonically decreasing and self-inverse.
- $b(x \wedge y)$ depends only on b(y) and b(x|y).
- Consistency
 - If a conclusion can be reasoned in more than one way, then every way should lead to the same answer.
 - Beliefs always take into account all relevant evidence.
 - Equivalent states of knowledge are represented by equivalent plausibility assignments.

Consequence: Belief functions (e.g. b(x), b(x|y), b(x, y)) must be isomorphic to probabilities, satisfying all the usual laws, including Bayes rule. (See Jaynes, *Probability Theory: The Logic of Science*)

• Probabilities are non-negative $P(x) \ge 0 \ \forall x$.

- Probabilities are non-negative $P(x) \ge 0 \ \forall x$.
- ▶ Probabilities normalise: $\sum_{x \in \mathcal{X}} P(x) = 1$ for distributions if *x* is a discrete variable and $\int_{-\infty}^{+\infty} p(x) dx = 1$ for probability densities over continuous variables

- Probabilities are non-negative $P(x) \ge 0 \ \forall x$.
- ▶ Probabilities normalise: $\sum_{x \in \mathcal{X}} P(x) = 1$ for distributions if *x* is a discrete variable and $\int_{-\infty}^{+\infty} p(x) dx = 1$ for probability densities over continuous variables
- The joint probability of x and y is: P(x, y).

- ▶ Probabilities are non-negative $P(x) \ge 0 \forall x$.
- ▶ Probabilities normalise: $\sum_{x \in \mathcal{X}} P(x) = 1$ for distributions if *x* is a discrete variable and $\int_{-\infty}^{+\infty} p(x) dx = 1$ for probability densities over continuous variables
- The joint probability of x and y is: P(x, y).
- The marginal probability of x is: $P(x) = \sum_{y} P(x, y)$, assuming y is discrete.

- ▶ Probabilities are non-negative $P(x) \ge 0 \forall x$.
- ▶ Probabilities normalise: $\sum_{x \in \mathcal{X}} P(x) = 1$ for distributions if *x* is a discrete variable and $\int_{-\infty}^{+\infty} p(x) dx = 1$ for probability densities over continuous variables
- The joint probability of x and y is: P(x, y).
- The marginal probability of x is: $P(x) = \sum_{y} P(x, y)$, assuming y is discrete.
- The conditional probability of x given y is: P(x|y) = P(x, y)/P(y)

- ▶ Probabilities are non-negative $P(x) \ge 0 \forall x$.
- ▶ Probabilities normalise: $\sum_{x \in \mathcal{X}} P(x) = 1$ for distributions if *x* is a discrete variable and $\int_{-\infty}^{+\infty} p(x) dx = 1$ for probability densities over continuous variables
- The joint probability of x and y is: P(x, y).
- The marginal probability of x is: $P(x) = \sum_{y} P(x, y)$, assuming y is discrete.
- The conditional probability of x given y is: P(x|y) = P(x,y)/P(y)
- Bayes Rule:

$$P(x,y) = P(x)P(y|x) = P(y)P(x|y) \implies P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

Assume you are willing to accept bets with odds proportional to the strength of your beliefs. That is, b(x) = 0.9 implies that you will accept a bet:

 $x \text{ at } 1:9 \Rightarrow \begin{cases} x & \text{is true} & \text{win} \geq \pounds 1 \\ x & \text{is false} & \text{lose} \quad \pounds 9 \end{cases}$

Then, unless your beliefs satisfy the rules of probability theory, including Bayes rule, there exists a set of simultaneous bets (called a "Dutch Book") which you are willing to accept, and for which **you are guaranteed to lose money, no matter what the outcome**. E.g. suppose $A \cap B = \emptyset$, then

Assume you are willing to accept bets with odds proportional to the strength of your beliefs. That is, b(x) = 0.9 implies that you will accept a bet:

 $x \text{ at } 1:9 \Rightarrow \begin{cases} x & \text{is true} & \text{win} \geq \pounds 1 \\ x & \text{is false} & \text{lose} \quad \pounds 9 \end{cases}$

Then, unless your beliefs satisfy the rules of probability theory, including Bayes rule, there exists a set of simultaneous bets (called a "Dutch Book") which you are willing to accept, and for which **you are guaranteed to lose money, no matter what the outcome**. E.g. suppose $A \cap B = \emptyset$, then

$$\left\{ \begin{array}{cc} b(A) &= 0.3\\ b(B) &= 0.2\\ b(A \cup B) &= 0.6 \end{array} \right\} \Rightarrow \text{accept the bets} \left\{ \begin{array}{cc} \neg A & \text{at } 3:7\\ \neg B & \text{at } 2:8\\ A \cup B & \text{at } 4:6 \end{array} \right\}$$

Assume you are willing to accept bets with odds proportional to the strength of your beliefs. That is, b(x) = 0.9 implies that you will accept a bet:

 $x \text{ at } 1:9 \Rightarrow \begin{cases} x & \text{is true} & \text{win} \geq \pounds 1 \\ x & \text{is false} & \text{lose} \quad \pounds 9 \end{cases}$

Then, unless your beliefs satisfy the rules of probability theory, including Bayes rule, there exists a set of simultaneous bets (called a "Dutch Book") which you are willing to accept, and for which **you are guaranteed to lose money, no matter what the outcome**. E.g. suppose $A \cap B = \emptyset$, then

$$\left\{\begin{array}{rrr} b(A) &= 0.3\\ b(B) &= 0.2\\ b(A \cup B) &= 0.6\end{array}\right\} \Rightarrow \text{accept the bets} \left\{\begin{array}{rrr} \neg A & \text{at } 3:7\\ \neg B & \text{at } 2:8\\ A \cup B & \text{at } 4:6\end{array}\right\}$$

But then:

 $\neg A \cap B \Rightarrow \min + 3 - 8 + 4 = -1$ $A \cap \neg B \Rightarrow \min - 7 + 2 + 4 = -1$ $\neg A \cap \neg B \Rightarrow \min + 3 + 2 - 6 = -1$

Assume you are willing to accept bets with odds proportional to the strength of your beliefs. That is, b(x) = 0.9 implies that you will accept a bet:

 $x \text{ at } 1:9 \Rightarrow \begin{cases} x & \text{is true} & \text{win} \geq \pounds 1 \\ x & \text{is false} & \text{lose} \quad \pounds 9 \end{cases}$

Then, unless your beliefs satisfy the rules of probability theory, including Bayes rule, there exists a set of simultaneous bets (called a "Dutch Book") which you are willing to accept, and for which **you are guaranteed to lose money, no matter what the outcome**. E.g. suppose $A \cap B = \emptyset$, then

$$\left\{\begin{array}{rrr} b(A) &= 0.3\\ b(B) &= 0.2\\ b(A \cup B) &= 0.6\end{array}\right\} \Rightarrow \text{accept the bets} \left\{\begin{array}{rrr} \neg A & \text{at } 3:7\\ \neg B & \text{at } 2:8\\ A \cup B & \text{at } 4:6\end{array}\right\}$$

But then:

 $\neg A \cap B \Rightarrow \min + 3 - 8 + 4 = -1$ $A \cap \neg B \Rightarrow \min - 7 + 2 + 4 = -1$ $\neg A \cap \neg B \Rightarrow \min + 3 + 2 - 6 = -1$

The only way to guard against Dutch Books is to ensure that your beliefs are coherent: i.e. satisfy the rules of probability.

Bayesian learning

Apply the basic rules of probability to learning from data.

Problem specification:

Data: $\mathcal{D} = \{x_1, \dots, x_n\}$ Models: $\mathcal{M}_1, \mathcal{M}_2$, etc. Parameters: θ_i (per model) Prior probability of models: $P(\mathcal{M}_i)$. Prior probabilities of model parameters: $P(\theta_i | \mathcal{M}_i)$ Model of data given parameters (likelihood model): $P(x | \theta_i, \mathcal{M}_i)$

Bayesian learning

Apply the basic rules of probability to learning from data.

Problem specification:

Data: $\mathcal{D} = \{x_1, \dots, x_n\}$ Models: $\mathcal{M}_1, \mathcal{M}_2$, etc. Parameters: θ_i (per model) Prior probability of models: $P(\mathcal{M}_i)$. Prior probabilities of model parameters: $P(\theta_i | \mathcal{M}_i)$ Model of data given parameters (likelihood model): $P(x | \theta_i, \mathcal{M}_i)$

Data probability (likelihood)

$$P(\mathcal{D}|\theta_i, \mathcal{M}_i) = \prod_{j=1}^n P(x_j|\theta_i, \mathcal{M}_i) \equiv \mathcal{L}(\theta_i)$$

(provided the data are independently and identically distributed (iid).
Bayesian learning

Apply the basic rules of probability to learning from data.

Problem specification:

Data: $\mathcal{D} = \{x_1, \dots, x_n\}$ Models: $\mathcal{M}_1, \mathcal{M}_2$, etc. Parameters: θ_i (per model) Prior probability of models: $P(\mathcal{M}_i)$. Prior probabilities of model parameters: $P(\theta_i | \mathcal{M}_i)$ Model of data given parameters (likelihood model): $P(x | \theta_i, \mathcal{M}_i)$

Data probability (likelihood)

$$P(\mathcal{D}|\theta_i, \mathcal{M}_i) = \prod_{j=1}^n P(x_j|\theta_i, \mathcal{M}_i) \equiv \mathcal{L}(\theta_i)$$

(provided the data are independently and identically distributed (iid).

Parameter learning (posterior):

$$\mathsf{P}(heta_i | \mathcal{D}, \mathcal{M}_i) = rac{\mathsf{P}(\mathcal{D} | heta_i, \mathcal{M}_i) \mathsf{P}(heta_i | \mathcal{M}_i)}{\mathsf{P}(\mathcal{D} | \mathcal{M}_i)}; \quad \mathsf{P}(\mathcal{D} | \mathcal{M}_i) = \int d heta_i \ \mathsf{P}(\mathcal{D} | heta_i, \mathcal{M}_i) \mathsf{P}(heta_i | \mathcal{M}_i)$$

Bayesian learning

Apply the basic rules of probability to learning from data.

Problem specification:

Data: $\mathcal{D} = \{x_1, \dots, x_n\}$ Models: $\mathcal{M}_1, \mathcal{M}_2$, etc. Parameters: θ_i (per model) Prior probability of models: $P(\mathcal{M}_i)$. Prior probabilities of model parameters: $P(\theta_i | \mathcal{M}_i)$ Model of data given parameters (likelihood model): $P(x | \theta_i, \mathcal{M}_i)$

Data probability (likelihood)

$$P(\mathcal{D}|\theta_i, \mathcal{M}_i) = \prod_{j=1}^n P(x_j|\theta_i, \mathcal{M}_i) \equiv \mathcal{L}(\theta_i)$$

(provided the data are independently and identically distributed (iid).

Parameter learning (posterior):

$${m P}(heta_i | \mathcal{D}, \mathcal{M}_i) = rac{{m P}(\mathcal{D} | heta_i, \mathcal{M}_i) {m P}(heta_i | \mathcal{M}_i)}{{m P}(\mathcal{D} | \mathcal{M}_i)}; \quad {m P}(\mathcal{D} | \mathcal{M}_i) = \int {m d} heta_i \ {m P}(\mathcal{D} | heta_i, \mathcal{M}_i) {m P}(heta_i | \mathcal{M}_i)$$

 $P(\mathcal{D}|\mathcal{M}_i)$ is called the marginal likelihood or evidence for \mathcal{M}_i . It is proportional to the posterior probability model \mathcal{M}_i being the one that generated the data.

Bayesian learning

Apply the basic rules of probability to learning from data.

Problem specification:

Data: $\mathcal{D} = \{x_1, \dots, x_n\}$ Models: $\mathcal{M}_1, \mathcal{M}_2$, etc. Parameters: θ_i (per model) Prior probability of models: $P(\mathcal{M}_i)$. Prior probabilities of model parameters: $P(\theta_i | \mathcal{M}_i)$ Model of data given parameters (likelihood model): $P(x | \theta_i, \mathcal{M}_i)$

Data probability (likelihood)

$$P(\mathcal{D}|\theta_i, \mathcal{M}_i) = \prod_{j=1}^n P(x_j|\theta_i, \mathcal{M}_i) \equiv \mathcal{L}(\theta_i)$$

(provided the data are independently and identically distributed (iid).

Parameter learning (posterior):

$$\mathsf{P}(heta_i | \mathcal{D}, \mathcal{M}_i) = rac{\mathsf{P}(\mathcal{D} | heta_i, \mathcal{M}_i) \mathsf{P}(heta_i | \mathcal{M}_i)}{\mathsf{P}(\mathcal{D} | \mathcal{M}_i)}; \quad \mathsf{P}(\mathcal{D} | \mathcal{M}_i) = \int \mathsf{d} heta_i \ \mathsf{P}(\mathcal{D} | heta_i, \mathcal{M}_i) \mathsf{P}(heta_i | \mathcal{M}_i)$$

 $P(\mathcal{D}|\mathcal{M}_i)$ is called the marginal likelihood or evidence for \mathcal{M}_i . It is proportional to the posterior probability model \mathcal{M}_i being the one that generated the data.

Model selection:

$$P(\mathcal{M}_i|\mathcal{D}) = rac{P(\mathcal{D}|\mathcal{M}_i)P(\mathcal{M}_i)}{P(\mathcal{D})}$$

Coin toss: One parameter q — the probability of obtaining *heads* So our space of models is the set of distributions over $q \in [0, 1]$.

Coin toss: One parameter q — the probability of obtaining *heads* So our space of models is the set of distributions over $q \in [0, 1]$. Learner A believes model \mathcal{M}_A : all values of q are equally plausible;



Coin toss: One parameter q — the probability of obtaining *heads* So our space of models is the set of distributions over $q \in [0, 1]$. Learner A believes model \mathcal{M}_A : all values of q are equally plausible; Learner B believes model \mathcal{M}_B : more plausible that the coin is "fair" ($q \approx 0.5$) than "biased".



Coin toss: One parameter q — the probability of obtaining *heads* So our space of models is the set of distributions over $q \in [0, 1]$. Learner A believes model \mathcal{M}_A : all values of q are equally plausible; Learner B believes model \mathcal{M}_B : more plausible that the coin is "fair" ($q \approx 0.5$) than "biased".



Both prior beliefs can be described by the Beta distribution:

$$p(q|lpha_1, lpha_2) = rac{q^{(lpha_1-1)}(1-q)^{(lpha_2-1)}}{B(lpha_1, lpha_2)} = ext{Beta}(q|lpha_1, lpha_2)$$

Coin toss: One parameter q — the probability of obtaining *heads* So our space of models is the set of distributions over $q \in [0, 1]$. Learner A believes model \mathcal{M}_A : all values of q are equally plausible; Learner B believes model \mathcal{M}_B : more plausible that the coin is "fair" ($q \approx 0.5$) than "biased".



Both prior beliefs can be described by the Beta distribution:

$$p(q|lpha_1, lpha_2) = rac{q^{(lpha_1-1)}(1-q)^{(lpha_2-1)}}{B(lpha_1, lpha_2)} = ext{Beta}(q|lpha_1, lpha_2)$$

Now we observe a toss. Two possible outcomes:

p(H|q) = q p(T|q) = 1 - q

Now we observe a toss. Two possible outcomes:

p(H|q) = q p(T|q) = 1 - q

Suppose our single coin toss comes out heads

Now we observe a toss. Two possible outcomes:

p(H|q) = q p(T|q) = 1 - q

Suppose our single coin toss comes out heads

The probability of the observed data (likelihood) is:

 $p(\mathsf{H}|q) = q$

Now we observe a toss. Two possible outcomes:

p(H|q) = q p(T|q) = 1 - q

Suppose our single coin toss comes out heads

The probability of the observed data (likelihood) is:

p(H|q) = q

Using Bayes Rule, we multiply the prior, p(q) by the likelihood and renormalise to get the posterior probability:

 $p(q|H) = \frac{p(q)p(H|q)}{p(H)}$

Now we observe a toss. Two possible outcomes:

p(H|q) = q p(T|q) = 1 - q

Suppose our single coin toss comes out heads

The probability of the observed data (likelihood) is:

p(H|q) = q

Using Bayes Rule, we multiply the prior, p(q) by the likelihood and renormalise to get the posterior probability:

$$p(q|\mathsf{H}) = rac{p(q)p(\mathsf{H}|q)}{p(\mathsf{H})} \propto q \operatorname{Beta}(q|lpha_1, lpha_2)$$

Now we observe a toss. Two possible outcomes:

p(H|q) = q p(T|q) = 1 - q

Suppose our single coin toss comes out heads

The probability of the observed data (likelihood) is:

 $p(\mathsf{H}|q) = q$

Using Bayes Rule, we multiply the prior, p(q) by the likelihood and renormalise to get the posterior probability:

$$p(q|H) = \frac{p(q)p(H|q)}{p(H)} \propto q \operatorname{Beta}(q|\alpha_1, \alpha_2)$$
$$\propto q q^{(\alpha_1-1)}(1-q)^{(\alpha_2-1)}$$

Now we observe a toss. Two possible outcomes:

p(H|q) = q p(T|q) = 1 - q

Suppose our single coin toss comes out heads

The probability of the observed data (likelihood) is:

p(H|q) = q

Using Bayes Rule, we multiply the prior, p(q) by the likelihood and renormalise to get the posterior probability:

$$p(q|H) = \frac{p(q)p(H|q)}{p(H)} \propto q \operatorname{Beta}(q|\alpha_1, \alpha_2)$$
$$\propto q q^{(\alpha_1 - 1)}(1 - q)^{(\alpha_2 - 1)} = \operatorname{Beta}(q|\alpha_1 + 1, \alpha_2)$$



What about multiple tosses?

What about multiple tosses? Suppose we observe $\mathcal{D} = \{ H H T H T T \}$:

$$p(\{HHTHTT\}|q) = qq(1-q)q(1-q)(1-q) = q^{3}(1-q)^{3}$$

What about multiple tosses? Suppose we observe $\mathcal{D} = \{ H H T H T T \}$:

$$p(\{ H H T H T T \}|q) = qq(1-q)q(1-q)(1-q) = q^{3}(1-q)^{3}$$

This is still straightforward:

What about multiple tosses? Suppose we observe $\mathcal{D} = \{ H H T H T T \}$:

$$p(\{ H H T H T T \}|q) = qq(1-q)q(1-q)(1-q) = q^{3}(1-q)^{3}$$

This is still straightforward:

$$p(q|\mathcal{D}) = \frac{p(q)p(\mathcal{D}|q)}{p(\mathcal{D})}$$

What about multiple tosses? Suppose we observe $\mathcal{D} = \{ H H T H T T \}$:

$$p(\{ H H T H T T \}|q) = qq(1-q)q(1-q)(1-q) = q^{3}(1-q)^{3}$$

This is still straightforward:

$$p(q|\mathcal{D}) = rac{p(q)p(\mathcal{D}|q)}{p(\mathcal{D})} \propto q^3(1-q)^3 \operatorname{Beta}(q|\alpha_1, \alpha_2)$$

What about multiple tosses? Suppose we observe $\mathcal{D} = \{ H H T H T T \}$:

$$p(\{ H H T H T T \}|q) = qq(1-q)q(1-q)(1-q) = q^3(1-q)^3$$

This is still straightforward:

$$p(q|\mathcal{D}) = rac{p(q)p(\mathcal{D}|q)}{p(\mathcal{D})} \propto q^3(1-q)^3 \operatorname{Beta}(q|\alpha_1,\alpha_2)$$

 \propto Beta($q|\alpha_1 + 3, \alpha_2 + 3$)

What about multiple tosses? Suppose we observe $\mathcal{D} = \{ H H T H T T \}$:

$$p(\{HHTHTT\}|q) = qq(1-q)q(1-q)(1-q) = q^{3}(1-q)^{3}$$

This is still straightforward:

$$p(q|\mathcal{D}) = rac{p(q)p(\mathcal{D}|q)}{p(\mathcal{D})} \propto q^3(1-q)^3 \operatorname{Beta}(q|lpha_1, lpha_2)$$

 \propto Beta($q|\alpha_1 + 3, \alpha_2 + 3$)



Updating the prior to form the posterior was particularly easy in these examples. This is because we used a conjugate prior for an exponential family likelihood.

Updating the prior to form the posterior was particularly easy in these examples. This is because we used a conjugate prior for an exponential family likelihood. Exponential family distributions take the form:

 $P(x|\theta) = g(\theta)f(x)e^{\phi(\theta)^{\mathsf{T}}\mathsf{T}(x)}$

with $g(\theta)$ the normalising constant.

Updating the prior to form the posterior was particularly easy in these examples. This is because we used a conjugate prior for an exponential family likelihood. Exponential family distributions take the form:

$$P(x|\theta) = g(\theta)f(x)e^{\phi(\theta)^{\mathsf{T}}\mathsf{T}(x)}$$

with $g(\theta)$ the normalising constant. Given *n* iid observations,

$$P(\{x_i\}|\theta) = \prod_i P(x_i|\theta) = g(\theta)^n e^{\phi(\theta)^T \left(\sum_i T(x_i)\right)} \prod_i f(x_i)$$

Updating the prior to form the posterior was particularly easy in these examples. This is because we used a conjugate prior for an exponential family likelihood. Exponential family distributions take the form:

$$P(x|\theta) = g(\theta)f(x)e^{\phi(\theta)^{\mathsf{T}}\mathsf{T}(x)}$$

with $g(\theta)$ the normalising constant. Given *n* iid observations,

$$P(\{x_i\}|\theta) = \prod_i P(x_i|\theta) = g(\theta)^n e^{\phi(\theta)^{\mathsf{T}} \left(\sum_i \mathsf{T}(x_i)\right)} \prod_i f(x_i)$$

Thus, if the prior takes the conjugate form

$$P(\theta) = F(\tau, \nu)g(\theta)^{\nu}e^{\phi(\theta)^{\mathsf{T}}\tau}$$

with $F(\tau, \nu)$ the normaliser

Updating the prior to form the posterior was particularly easy in these examples. This is because we used a conjugate prior for an exponential family likelihood. Exponential family distributions take the form:

$$P(x|\theta) = g(\theta)f(x)e^{\phi(\theta)^{\mathsf{T}}\mathsf{T}(x)}$$

with $g(\theta)$ the normalising constant. Given *n* iid observations,

$$P(\{x_i\}|\theta) = \prod_i P(x_i|\theta) = g(\theta)^n e^{\phi(\theta)^{\mathsf{T}} \left(\sum_i \mathsf{T}(x_i)\right)} \prod_i f(x_i)$$

Thus, if the prior takes the conjugate form

$$P(\theta) = F(\tau, \nu)g(\theta)^{\nu}e^{\phi(\theta)^{\mathsf{T}}\tau}$$

with $F(\tau, \nu)$ the normaliser, then the posterior is

$$\mathcal{P}(\theta|\{x_i\}) \propto \mathcal{P}(\{x_i\}|\theta)\mathcal{P}(\theta) \propto g(\theta)^{\nu+n} e^{\phi(\theta)^{\mathsf{T}} \left(\tau + \sum_i \mathsf{T}(x_i)\right)}$$

with the normaliser given by $F(\tau + \sum_{i} \mathbf{T}(x_i), \nu + n)$.

The posterior given an exponential family likelihood and conjugate prior is:

$$P(\theta|\{x_i\}) = F(\tau + \sum_i \mathbf{T}(x_i), \nu + n)g(\theta)^{\nu+n} \exp\left[\phi(\theta)^{\mathsf{T}} \left(\tau + \sum_i \mathbf{T}(x_i)\right)\right]$$

Here,

- $\phi(\theta)$ is the vector of natural parameters
- $\sum_{i} \mathbf{T}(x_i)$ is the vector of sufficient statistics
 - au are pseudo-observations which define the prior
 - ν is the scale of the prior (need not be an integer)

As new data come in, each one increments the sufficient statistics vector and the scale to define the posterior.

The posterior given an exponential family likelihood and conjugate prior is:

$$P(\theta|\{x_i\}) = F(\tau + \sum_i \mathbf{T}(x_i), \nu + n)g(\theta)^{\nu + n} \exp\left[\phi(\theta)^{\mathsf{T}} \left(\tau + \sum_i \mathbf{T}(x_i)\right)\right]$$

Here,

- $\phi(heta)$ is the vector of natural parameters
- $\sum_{i} \mathbf{T}(x_i)$ is the vector of sufficient statistics
 - au are pseudo-observations which define the prior
 - ν is the scale of the prior (need not be an integer)

As new data come in, each one increments the sufficient statistics vector and the scale to define the posterior.

The prior appears to be based on "pseudo-observations", but:

- 1. This is different to applying Bayes' rule. No prior! Sometimes we can take a uniform prior (say on [0, 1] for *q*), but for unbounded *θ*, there may be no equivalent.
- 2. A valid conjugate prior might have non-integral ν or impossible τ , with no likelihood equivalent.

Distributions are not always written in their natural exponential form.

Distributions are not always written in their natural exponential form.

The Bernoulli distribution (a single coin flip) with parameter q and observation $x \in \{0, 1\}$, can be written:

$$P(x|q) = q^{x}(1-q)^{(1-x)}$$

= $e^{x \log q + (1-x) \log(1-q)}$
= $e^{\log(1-q) + x \log(q/(1-q))}$
= $(1-q)e^{\log(q/(1-q))x}$

Distributions are not always written in their natural exponential form.

The Bernoulli distribution (a single coin flip) with parameter q and observation $x \in \{0, 1\}$, can be written:

$$P(x|q) = q^{x} (1-q)^{(1-x)}$$

= $e^{x \log q + (1-x) \log(1-q)}$
= $e^{\log(1-q) + x \log(q/(1-q))}$
= $(1-q)e^{\log(q/(1-q))x}$

So the natural parameter is the log odds $\log(q/(1-q))$, and the sufficient stats (for multiple tosses) is the number of heads.

Distributions are not always written in their natural exponential form.

The Bernoulli distribution (a single coin flip) with parameter q and observation $x \in \{0, 1\}$, can be written:

$$P(x|q) = q^{x} (1-q)^{(1-x)}$$

= $e^{x \log q + (1-x) \log(1-q)}$
= $e^{\log(1-q) + x \log(q/(1-q))}$
= $(1-q)e^{\log(q/(1-q))x}$

So the natural parameter is the log odds $\log(q/(1-q))$, and the sufficient stats (for multiple tosses) is the number of heads.

The conjugate prior is

$$P(q) = F(\tau, \nu) (1 - q)^{\nu} e^{\log(q/(1-q))\tau}$$

= $F(\tau, \nu) (1 - q)^{\nu} e^{\tau \log q - \tau \log(1-q)}$
= $F(\tau, \nu) (1 - q)^{\nu - \tau} q^{\tau}$

which has the form of the Beta distribution $\Rightarrow F(\tau, \nu) = 1/B(\tau + 1, \nu - \tau + 1)$.

Distributions are not always written in their natural exponential form.

The Bernoulli distribution (a single coin flip) with parameter q and observation $x \in \{0, 1\}$, can be written:

$$P(x|q) = q^{x} (1-q)^{(1-x)}$$

= $e^{x \log q + (1-x) \log(1-q)}$
= $e^{\log(1-q) + x \log(q/(1-q))}$
= $(1-q)e^{\log(q/(1-q))x}$

So the natural parameter is the log odds $\log(q/(1-q))$, and the sufficient stats (for multiple tosses) is the number of heads.

The conjugate prior is

$$\begin{split} P(q) &= F(\tau,\nu) \; (1-q)^{\nu} e^{\log(q/(1-q))\tau} \\ &= F(\tau,\nu) \; (1-q)^{\nu} e^{\tau \log q - \tau \log(1-q)} \\ &= F(\tau,\nu) \; (1-q)^{\nu-\tau} q^{\tau} \end{split}$$

which has the form of the Beta distribution $\Rightarrow F(\tau, \nu) = 1/B(\tau + 1, \nu - \tau + 1)$. In general, then, the posterior will be $P(q|\{x_i\}) = \text{Beta}(\alpha_1, \alpha_2)$, with

$$\alpha_1 = 1 + \tau + \sum_i x_i \qquad \qquad \alpha_2 = 1 + (\nu + n) - \left(\tau + \sum_i x_i\right)$$

Distributions are not always written in their natural exponential form.

The Bernoulli distribution (a single coin flip) with parameter q and observation $x \in \{0, 1\}$, can be written:

$$P(x|q) = q^{x} (1-q)^{(1-x)}$$

= $e^{x \log q + (1-x) \log(1-q)}$
= $e^{\log(1-q) + x \log(q/(1-q))}$
= $(1-q)e^{\log(q/(1-q))x}$

So the natural parameter is the log odds $\log(q/(1-q))$, and the sufficient stats (for multiple tosses) is the number of heads.

The conjugate prior is

$$\begin{split} P(q) &= F(\tau,\nu) \; (1-q)^{\nu} e^{\log(q/(1-q))\tau} \\ &= F(\tau,\nu) \; (1-q)^{\nu} e^{\tau \log q - \tau \log(1-q)} \\ &= F(\tau,\nu) \; (1-q)^{\nu-\tau} q^{\tau} \end{split}$$

which has the form of the Beta distribution $\Rightarrow F(\tau, \nu) = 1/B(\tau + 1, \nu - \tau + 1)$. In general, then, the posterior will be $P(q|\{x_i\}) = \text{Beta}(\alpha_1, \alpha_2)$, with

$$\alpha_1 = 1 + \tau + \sum_i x_i \qquad \qquad \alpha_2 = 1 + (\nu + n) - \left(\tau + \sum_i x_i\right)$$

If we observe a head, we add 1 to the sufficient statistic $\sum x_i$, and also 1 to the count *n*. This increments α_1 .
Conjugacy in the coin flip

Distributions are not always written in their natural exponential form.

The Bernoulli distribution (a single coin flip) with parameter q and observation $x \in \{0, 1\}$, can be written:

$$P(x|q) = q^{x}(1-q)^{(1-x)}$$

= $e^{x \log q + (1-x) \log(1-q)}$
= $e^{\log(1-q) + x \log(q/(1-q))}$
= $(1-q)e^{\log(q/(1-q))x}$

So the natural parameter is the log odds $\log(q/(1 - q))$, and the sufficient stats (for multiple tosses) is the number of heads.

The conjugate prior is

$$\begin{split} P(q) &= F(\tau,\nu) \; (1-q)^{\nu} e^{\log(q/(1-q))\tau} \\ &= F(\tau,\nu) \; (1-q)^{\nu} e^{\tau \log q - \tau \log(1-q)} \\ &= F(\tau,\nu) \; (1-q)^{\nu-\tau} q^{\tau} \end{split}$$

which has the form of the Beta distribution $\Rightarrow F(\tau, \nu) = 1/B(\tau + 1, \nu - \tau + 1)$. In general, then, the posterior will be $P(q|\{x_i\}) = \text{Beta}(\alpha_1, \alpha_2)$, with

$$\alpha_1 = 1 + \tau + \sum_i X_i \qquad \qquad \alpha_2 = 1 + (\nu + n) - \left(\tau + \sum_i X_i\right)$$

If we observe a head, we add 1 to the sufficient statistic $\sum x_i$, and also 1 to the count *n*. This increments α_1 . If we observe a tail we add 1 to *n*, but not to $\sum x_i$, incrementing α_2 .

We have seen how to update posteriors within each model. To study the choice of model, consider two more extreme models: "fair" and "bent".

We have seen how to update posteriors within each model. To study the choice of model, consider two more extreme models: "fair" and "bent". A priori, we may think that "fair" is more probable, eg:

 $p(fair) = 0.8, \quad p(bent) = 0.2$

We have seen how to update posteriors within each model. To study the choice of model, consider two more extreme models: "fair" and "bent". A priori, we may think that "fair" is more probable, eg:

 $p(fair) = 0.8, \quad p(bent) = 0.2$

For the bent coin, we assume all parameter values are equally likely, whilst the fair coin has a fixed probability:



We have seen how to update posteriors within each model. To study the choice of model, consider two more extreme models: "fair" and "bent". A priori, we may think that "fair" is more probable, eg:

 $p(fair) = 0.8, \quad p(bent) = 0.2$

For the bent coin, we assume all parameter values are equally likely, whilst the fair coin has a fixed probability:



We make 10 tosses, and get: $\mathcal{D} = (T H T H T T T T T)$.

Which model should we prefer a posteriori (i.e. after seeing the data)?

Which model should we prefer a posteriori (i.e. after seeing the data)?

The evidence for the fair model is:

 $P(D|\text{fair}) = (1/2)^{10} \approx 0.001$

Which model should we prefer a posteriori (i.e. after seeing the data)?

The evidence for the fair model is:

 $P(D|\text{fair}) = (1/2)^{10} \approx 0.001$

and for the bent model is:

$$P(\mathcal{D}|\mathsf{bent}) = \int dq \ P(\mathcal{D}|q,\mathsf{bent})p(q|\mathsf{bent})$$

Which model should we prefer a posteriori (i.e. after seeing the data)?

The evidence for the fair model is:

$$P(D|\text{fair}) = (1/2)^{10} \approx 0.001$$

and for the bent model is:

$$P(\mathcal{D}|\text{bent}) = \int dq \ P(\mathcal{D}|q, \text{bent}) p(q|\text{bent}) = \int dq \ q^2(1-q)^8$$

Which model should we prefer a posteriori (i.e. after seeing the data)?

The evidence for the fair model is:

$$P(D|\text{fair}) = (1/2)^{10} \approx 0.001$$

and for the bent model is:

$$\mathcal{P}(\mathcal{D}|\mathsf{bent}) = \int dq \; \mathcal{P}(\mathcal{D}|q,\mathsf{bent}) \mathcal{p}(q|\mathsf{bent}) \; = \int dq \; q^2 (1-q)^8 \; = \mathrm{B}(3,9) pprox 0.002$$

Which model should we prefer a posteriori (i.e. after seeing the data)?

The evidence for the fair model is:

 $P(D|\text{fair}) = (1/2)^{10} \approx 0.001$

and for the bent model is:

$$P(\mathcal{D}|\text{bent}) = \int dq \ P(\mathcal{D}|q,\text{bent})p(q|\text{bent}) = \int dq \ q^2(1-q)^8 = B(3,9) \approx 0.002$$

Thus, the posterior for the models, by Bayes rule:

 $P(\text{fair}|\mathcal{D}) \propto 0.0008, \qquad P(\text{bent}|\mathcal{D}) \propto 0.0004,$

ie, a two-thirds probability that the coin is fair.

Which model should we prefer a posteriori (i.e. after seeing the data)?

The evidence for the fair model is:

 $P(D|\text{fair}) = (1/2)^{10} \approx 0.001$

and for the bent model is:

$$P(\mathcal{D}|\text{bent}) = \int dq \ P(\mathcal{D}|q,\text{bent})p(q|\text{bent}) = \int dq \ q^2(1-q)^8 = B(3,9) \approx 0.002$$

Thus, the posterior for the models, by Bayes rule:

 $P(\text{fair}|\mathcal{D}) \propto 0.0008, \qquad P(\text{bent}|\mathcal{D}) \propto 0.0004,$

ie, a two-thirds probability that the coin is fair.

How do we make predictions?

Which model should we prefer a posteriori (i.e. after seeing the data)?

The evidence for the fair model is:

 $P(D|\text{fair}) = (1/2)^{10} \approx 0.001$

and for the bent model is:

$$P(\mathcal{D}|\text{bent}) = \int dq \ P(\mathcal{D}|q,\text{bent})p(q|\text{bent}) = \int dq \ q^2(1-q)^8 = B(3,9) \approx 0.002$$

Thus, the posterior for the models, by Bayes rule:

 $P(\text{fair}|\mathcal{D}) \propto 0.0008, \qquad P(\text{bent}|\mathcal{D}) \propto 0.0004,$

ie, a two-thirds probability that the coin is fair.

How do we make predictions? Could choose the fair model (model selection).

Which model should we prefer a posteriori (i.e. after seeing the data)?

The evidence for the fair model is:

 $P(D|\text{fair}) = (1/2)^{10} \approx 0.001$

and for the bent model is:

$$P(\mathcal{D}|\text{bent}) = \int dq \ P(\mathcal{D}|q,\text{bent})p(q|\text{bent}) = \int dq \ q^2(1-q)^8 = B(3,9) \approx 0.002$$

Thus, the posterior for the models, by Bayes rule:

 $P(\text{fair}|\mathcal{D}) \propto 0.0008, \qquad P(\text{bent}|\mathcal{D}) \propto 0.0004,$

ie, a two-thirds probability that the coin is fair.

How do we make predictions? Could choose the fair model (model selection). Or could weight the predictions from each model by their probability (model averaging). Probability of H at next toss is:

$$P(\mathsf{H}|\mathcal{D}) = P(\mathsf{H}|\mathcal{D},\mathsf{fair})P(\mathsf{fair}|\mathcal{D}) + P(\mathsf{H}|\mathcal{D},\mathsf{bent})P(\mathsf{bent}|\mathcal{D}) = \frac{2}{3} \times \frac{1}{2} + \frac{1}{3} \times \frac{3}{12} = \frac{5}{12}.$$

The Bayesian probabilistic prescription tells us how to reason about models and their parameters.

The Bayesian probabilistic prescription tells us how to reason about models and their parameters. But it is often impractical for realistic models (outside the exponential family).

The Bayesian probabilistic prescription tells us how to reason about models and their parameters. But it is often impractical for realistic models (outside the exponential family).

Point estimates of parameters or other predictions

The Bayesian probabilistic prescription tells us how to reason about models and their parameters. But it is often impractical for realistic models (outside the exponential family).

- Point estimates of parameters or other predictions
 - Compute posterior and find single parameter that minimises expected loss.

$$\theta^{\mathsf{BP}} = \underset{\hat{\theta}}{\operatorname{argmin}} \left\langle L(\hat{\theta}, \theta) \right\rangle_{P(\theta \mid \mathcal{D})}$$

The Bayesian probabilistic prescription tells us how to reason about models and their parameters. But it is often impractical for realistic models (outside the exponential family).

- Point estimates of parameters or other predictions
 - Compute posterior and find single parameter that minimises expected loss.

$$\theta^{\mathsf{BP}} = \operatorname*{argmin}_{\hat{\theta}} \left\langle L(\hat{\theta}, \theta) \right\rangle_{P(\theta \mid \mathcal{D})}$$

• $\langle \theta \rangle_{P(\theta \mid D)}$ minimises squared loss.

The Bayesian probabilistic prescription tells us how to reason about models and their parameters. But it is often impractical for realistic models (outside the exponential family).

- Point estimates of parameters or other predictions
 - Compute posterior and find single parameter that minimises expected loss.

$$heta^{\mathsf{BP}} = \operatorname*{argmin}_{\hat{ heta}} \left\langle L(\hat{ heta}, heta) \right\rangle_{P(heta \mid \mathcal{D})}$$

• $\langle \theta \rangle_{P(\theta \mid D)}$ minimises squared loss.

Maximum a Posteriori (MAP) estimate: Assume a prior over the model parameters P(θ), and compute parameters that are most probable under the posterior:

 $\theta^{MAP} = \operatorname{argmax} P(\theta | D) = \operatorname{argmax} P(\theta) P(D | \theta).$

The Bayesian probabilistic prescription tells us how to reason about models and their parameters. But it is often impractical for realistic models (outside the exponential family).

- Point estimates of parameters or other predictions
 - Compute posterior and find single parameter that minimises expected loss.

$$heta^{\mathsf{BP}} = \operatorname*{argmin}_{\hat{ heta}} \left\langle L(\hat{ heta}, heta) \right\rangle_{P(heta \mid \mathcal{D})}$$

• $\langle \theta \rangle_{P(\theta \mid D)}$ minimises squared loss.

Maximum a Posteriori (MAP) estimate: Assume a prior over the model parameters P(θ), and compute parameters that are most probable under the posterior:

$$\theta^{MAP} = \operatorname{argmax} P(\theta | \mathcal{D}) = \operatorname{argmax} P(\theta) P(\mathcal{D} | \theta)$$
.

Equivalent to minimising the 0/1 loss.

The Bayesian probabilistic prescription tells us how to reason about models and their parameters. But it is often impractical for realistic models (outside the exponential family).

- Point estimates of parameters or other predictions
 - Compute posterior and find single parameter that minimises expected loss.

$$heta^{\mathsf{BP}} = \operatorname*{argmin}_{\hat{ heta}} \left\langle L(\hat{ heta}, heta)
ight
angle_{P(heta \mid \mathcal{D})}$$

• $\langle \theta \rangle_{P(\theta \mid D)}$ minimises squared loss.

Maximum a Posteriori (MAP) estimate: Assume a prior over the model parameters P(θ), and compute parameters that are most probable under the posterior:

$$\theta^{MAP} = \operatorname{argmax} P(\theta | \mathcal{D}) = \operatorname{argmax} P(\theta) P(\mathcal{D} | \theta)$$
.

- Equivalent to minimising the 0/1 loss.
- Maximum Likelihood (ML) Learning: No prior over the parameters. Compute parameter value that maximises the likelihood function alone:

The Bayesian probabilistic prescription tells us how to reason about models and their parameters. But it is often impractical for realistic models (outside the exponential family).

- Point estimates of parameters or other predictions
 - Compute posterior and find single parameter that minimises expected loss.

$$heta^{\mathsf{BP}} = \operatorname*{argmin}_{\hat{ heta}} \left\langle L(\hat{ heta}, heta)
ight
angle_{P(heta \mid \mathcal{D})}$$

• $\langle \theta \rangle_{P(\theta \mid D)}$ minimises squared loss.

Maximum a Posteriori (MAP) estimate: Assume a prior over the model parameters P(θ), and compute parameters that are most probable under the posterior:

 $\theta^{MAP} = \operatorname{argmax} P(\theta|\mathcal{D}) = \operatorname{argmax} P(\theta)P(\mathcal{D}|\theta).$

- Equivalent to minimising the 0/1 loss.
- Maximum Likelihood (ML) Learning: No prior over the parameters. Compute parameter value that maximises the likelihood function alone:

 $\theta^{\mathsf{ML}} = \operatorname{argmax} P(\mathcal{D}|\theta)$.

Parameterisation-independent.

The Bayesian probabilistic prescription tells us how to reason about models and their parameters. But it is often impractical for realistic models (outside the exponential family).

- Point estimates of parameters or other predictions
 - Compute posterior and find single parameter that minimises expected loss.

$$heta^{\mathsf{BP}} = \operatorname*{argmin}_{\hat{ heta}} \left\langle L(\hat{ heta}, heta)
ight
angle_{P(heta \mid \mathcal{D})}$$

• $\langle \theta \rangle_{P(\theta \mid D)}$ minimises squared loss.

• Maximum a Posteriori (MAP) estimate: Assume a prior over the model parameters $P(\theta)$, and compute parameters that are most probable under the posterior:

 $\theta^{MAP} = \operatorname{argmax} P(\theta|\mathcal{D}) = \operatorname{argmax} P(\theta)P(\mathcal{D}|\theta).$

- Equivalent to minimising the 0/1 loss.
- Maximum Likelihood (ML) Learning: No prior over the parameters. Compute parameter value that maximises the likelihood function alone:

- Parameterisation-independent.
- Approximations may allow us to recover samples from posterior, or to find a distribution which is close in some sense.

The Bayesian probabilistic prescription tells us how to reason about models and their parameters. But it is often impractical for realistic models (outside the exponential family).

- Point estimates of parameters or other predictions
 - Compute posterior and find single parameter that minimises expected loss.

$$heta^{\mathsf{BP}} = \operatorname*{argmin}_{\hat{ heta}} \left\langle L(\hat{ heta}, heta)
ight
angle_{P(heta \mid \mathcal{D})}$$

• $\langle \theta \rangle_{P(\theta \mid D)}$ minimises squared loss.

• Maximum a Posteriori (MAP) estimate: Assume a prior over the model parameters $P(\theta)$, and compute parameters that are most probable under the posterior:

 $\theta^{MAP} = \operatorname{argmax} P(\theta | \mathcal{D}) = \operatorname{argmax} P(\theta) P(\mathcal{D} | \theta).$

- Equivalent to minimising the 0/1 loss.
- Maximum Likelihood (ML) Learning: No prior over the parameters. Compute parameter value that maximises the likelihood function alone:

- Parameterisation-independent.
- Approximations may allow us to recover samples from posterior, or to find a distribution which is close in some sense.
- Choosing between these and other alternatives may be a matter of definition, of goals (loss function), or of practicality.

The Bayesian probabilistic prescription tells us how to reason about models and their parameters. But it is often impractical for realistic models (outside the exponential family).

- Point estimates of parameters or other predictions
 - Compute posterior and find single parameter that minimises expected loss.

$$heta^{\mathsf{BP}} = \operatorname*{argmin}_{\hat{ heta}} \left\langle L(\hat{ heta}, heta) \right\rangle_{P(\theta \mid \mathcal{D})}$$

• $\langle \theta \rangle_{P(\theta \mid D)}$ minimises squared loss.

• Maximum a Posteriori (MAP) estimate: Assume a prior over the model parameters $P(\theta)$, and compute parameters that are most probable under the posterior:

 $\theta^{MAP} = \operatorname{argmax} P(\theta|\mathcal{D}) = \operatorname{argmax} P(\theta)P(\mathcal{D}|\theta).$

- Equivalent to minimising the 0/1 loss.
- Maximum Likelihood (ML) Learning: No prior over the parameters. Compute parameter value that maximises the likelihood function alone:

- Parameterisation-independent.
- Approximations may allow us to recover samples from posterior, or to find a distribution which is close in some sense.
- Choosing between these and other alternatives may be a matter of definition, of goals (loss function), or of practicality.
- For the next few weeks we will look at ML and MAP learning in more complex models. We will then return to the fully Bayesian formulation for the few intersting cases where it is tractable. Approximations will be addressed in the second half of the course.

Modelling associations between variables



- Data set $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- with each data point a vector of *D* features: $\mathbf{x}_i = [x_{i1} \dots x_{iD}]$
- Assume data are i.i.d. (independent and identically distributed).

Modelling associations between variables



- Data set $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- with each data point a vector of *D* features: $\mathbf{x}_i = [x_{i1} \dots x_{iD}]$
- Assume data are i.i.d. (independent and identically distributed).

A simple forms of unsupervised (structure) learning: model the **mean** of the data and the **correlations** between the *D* features in the data.

Modelling associations between variables



- Data set $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- with each data point a vector of *D* features: $\mathbf{x}_i = [x_{i1} \dots x_{iD}]$
- Assume data are i.i.d. (independent and identically distributed).

A simple forms of unsupervised (structure) learning: model the **mean** of the data and the **correlations** between the *D* features in the data.

We can use a multivariate Gaussian model:

$$p(\mathbf{x}|\boldsymbol{\mu},\boldsymbol{\Sigma}) = \mathcal{N}(\boldsymbol{\mu},\boldsymbol{\Sigma}) = |2\pi\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^{\mathsf{T}}\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\right\}$$

Data set
$$\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$$
, likelihood: $p(\mathcal{D}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{n=1}^N p(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma})$

Goal: find μ and Σ that maximise likelihood

$$\mathcal{L} = \prod_{n=1}^{N} p(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Data set
$$\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$$
, likelihood: $p(\mathcal{D}|\mu, \Sigma) = \prod_{n=1}^N p(\mathbf{x}_n | \mu, \Sigma)$

Goal: find μ and Σ that maximise likelihood \Leftrightarrow maximise log likelihood:

$$\ell = \log \prod_{n=1}^{N} \rho(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Data set
$$\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$$
, likelihood: $p(\mathcal{D}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{n=1}^N p(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma})$

Goal: find μ and Σ that maximise likelihood \Leftrightarrow maximise log likelihood:

$$\ell = \log \prod_{n=1}^{N} p(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_n \log p(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Data set
$$\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$$
, likelihood: $p(\mathcal{D}|\mu, \Sigma) = \prod_{n=1}^N p(\mathbf{x}_n | \mu, \Sigma)$

Goal: find μ and Σ that maximise likelihood \Leftrightarrow maximise log likelihood:

$$\ell = \log \prod_{n=1}^{N} p(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n} \log p(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$
$$= -\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n} (\mathbf{x}_n - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$

Data set
$$\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$$
, likelihood: $p(\mathcal{D}|\mu, \Sigma) = \prod_{n=1}^N p(\mathbf{x}_n | \mu, \Sigma)$

Goal: find μ and Σ that maximise likelihood \Leftrightarrow maximise log likelihood:

$$\ell = \log \prod_{n=1}^{N} p(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n} \log p(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$
$$= -\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n} (\mathbf{x}_n - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$

Note: equivalently, minimise $-\ell$, which is *quadratic* in μ

Data set
$$\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$$
, likelihood: $\rho(\mathcal{D}|\mu, \Sigma) = \prod_{n=1}^N \rho(\mathbf{x}_n | \mu, \Sigma)$

Goal: find μ and Σ that maximise likelihood \Leftrightarrow maximise log likelihood:

$$\ell = \log \prod_{n=1}^{N} p(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n} \log p(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$
$$= -\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n} (\mathbf{x}_n - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$

Note: equivalently, minimise $-\ell$, which is *quadratic* in μ

Procedure: take derivatives and set to zero:

$$\frac{\partial \ell}{\partial \mu} = 0 \qquad \Rightarrow \qquad \hat{\mu} = \frac{1}{N} \sum_{n} \mathbf{x}_{n} \qquad \text{(sample mean)}$$
$$\frac{\partial \ell}{\partial \Sigma} = 0 \qquad \Rightarrow \qquad \hat{\Sigma} = \frac{1}{N} \sum_{n} (\mathbf{x}_{n} - \hat{\mu}) (\mathbf{x}_{n} - \hat{\mu})^{\mathsf{T}} \qquad \text{(sample covariance)}$$

Refresher – matrix derivatives of scalar forms

We will use the following facts:
$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \left[\mathbf{x}^{\mathsf{T}} A \mathbf{y} \right]$$
 (scalars equal their own transpose and trace)

We will use the following facts:

 $\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \mathsf{Tr} \left[\mathbf{x}^{\mathsf{T}} A \mathbf{y} \right]$ (scalars equal their own transpose and trace) Tr $[A] = \mathsf{Tr} \left[A^{\mathsf{T}} \right]$

We will use the following facts:

 $\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \begin{bmatrix} \mathbf{x}^{\mathsf{T}} A \mathbf{y} \end{bmatrix} \text{ (scalars equal their own transpose and trace)}$ $\operatorname{Tr} [A] = \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} \end{bmatrix} \qquad \operatorname{Tr} [ABC] = \operatorname{Tr} [CAB] = \operatorname{Tr} [BCA]$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \begin{bmatrix} \mathbf{x}^{\mathsf{T}} A \mathbf{y} \end{bmatrix} \text{ (scalars equal their own transpose and trace)}$$
$$\operatorname{Tr} [A] = \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} \end{bmatrix} \qquad \operatorname{Tr} [ABC] = \operatorname{Tr} [CAB] = \operatorname{Tr} [BCA]$$
$$\frac{\partial}{\partial A_{ij}} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix}$$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \begin{bmatrix} \mathbf{x}^{\mathsf{T}} A \mathbf{y} \end{bmatrix} \text{ (scalars equal their own transpose and trace)}$$
$$\operatorname{Tr} [A] = \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} \end{bmatrix} \qquad \operatorname{Tr} [ABC] = \operatorname{Tr} [CAB] = \operatorname{Tr} [BCA]$$
$$\frac{\partial}{\partial A_{ij}} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A_{ij}} \sum_{n} [A^{\mathsf{T}} B]_{nn}$$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \left[\mathbf{x}^{\mathsf{T}} A \mathbf{y} \right] \text{ (scalars equal their own transpose and trace)}$$
$$\operatorname{Tr} \left[A \right] = \operatorname{Tr} \left[A^{\mathsf{T}} \right] \qquad \operatorname{Tr} \left[A B C \right] = \operatorname{Tr} \left[C A B \right] = \operatorname{Tr} \left[B C A \right]$$
$$\frac{\partial}{\partial A_{ij}} \operatorname{Tr} \left[A^{\mathsf{T}} B \right] = \frac{\partial}{\partial A_{ij}} \sum_{n} \sum_{m} A_{nm}^{\mathsf{T}} B_{mn}$$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \begin{bmatrix} \mathbf{x}^{\mathsf{T}} A \mathbf{y} \end{bmatrix} \text{ (scalars equal their own transpose and trace)}$$
$$\operatorname{Tr} [A] = \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} \end{bmatrix} \qquad \operatorname{Tr} [ABC] = \operatorname{Tr} [CAB] = \operatorname{Tr} [BCA]$$
$$\frac{\partial}{\partial A_{ij}} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A_{ij}} \sum_{mn} A_{mn} B_{mn}$$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \left[\mathbf{x}^{\mathsf{T}} A \mathbf{y} \right] \text{ (scalars equal their own transpose and trace)}$$
$$\operatorname{Tr} \left[A \right] = \operatorname{Tr} \left[A^{\mathsf{T}} \right] \qquad \operatorname{Tr} \left[A B C \right] = \operatorname{Tr} \left[C A B \right] = \operatorname{Tr} \left[B C A \right]$$
$$\frac{\partial}{\partial A_{ij}} \operatorname{Tr} \left[A^{\mathsf{T}} B \right] = \frac{\partial}{\partial A_{ij}} \sum_{mn} A_{mn} B_{mn} = B_{ij}$$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \mathsf{Tr} \begin{bmatrix} \mathbf{x}^{\mathsf{T}} A \mathbf{y} \end{bmatrix} \text{ (scalars equal their own transpose and trace)}$$
$$\mathsf{Tr} [A] = \mathsf{Tr} \begin{bmatrix} A^{\mathsf{T}} \end{bmatrix} \qquad \mathsf{Tr} [ABC] = \mathsf{Tr} [CAB] = \mathsf{Tr} [BCA]$$
$$\frac{\partial}{\partial A_{ij}} \mathsf{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A_{ij}} \sum_{mn} A_{mn} B_{mn} = B_{ij}$$
$$\Rightarrow \frac{\partial}{\partial A} \mathsf{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = B$$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \mathsf{Tr} \begin{bmatrix} \mathbf{x}^{\mathsf{T}} A \mathbf{y} \end{bmatrix} \text{ (scalars equal their own transpose and trace)}$$
$$\mathsf{Tr} [A] = \mathsf{Tr} \begin{bmatrix} A^{\mathsf{T}} \end{bmatrix} \qquad \mathsf{Tr} [ABC] = \mathsf{Tr} [CAB] = \mathsf{Tr} [BCA]$$
$$\frac{\partial}{\partial A_{ij}} \mathsf{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A_{ij}} \sum_{mn} A_{mn} B_{mn} = B_{ij}$$
$$\Rightarrow \frac{\partial}{\partial A} \mathsf{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = B$$
$$\frac{\partial}{\partial A} \mathsf{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = B$$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \begin{bmatrix} \mathbf{x}^{\mathsf{T}} A \mathbf{y} \end{bmatrix} \text{ (scalars equal their own transpose and trace)}$$
$$\operatorname{Tr} [A] = \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} \end{bmatrix} \qquad \operatorname{Tr} [ABC] = \operatorname{Tr} [CAB] = \operatorname{Tr} [BCA]$$
$$\frac{\partial}{\partial A_{ij}} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A_{ij}} \sum_{mn} A_{mn} B_{mn} = B_{ij}$$
$$\Rightarrow \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = B$$
$$\frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = B$$
$$\frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B A C \end{bmatrix} = \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} F_1(A)^{\mathsf{T}} B F_2(A) C \end{bmatrix} \text{ with } F_1 \text{ and } F_2 \text{ both identity maps}$$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \begin{bmatrix} \mathbf{x}^{\mathsf{T}} A \mathbf{y} \end{bmatrix} \text{ (scalars equal their own transpose and trace)}$$
$$\operatorname{Tr} [A] = \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} \end{bmatrix} \qquad \operatorname{Tr} [ABC] = \operatorname{Tr} [CAB] = \operatorname{Tr} [BCA]$$
$$\frac{\partial}{\partial A_{ij}} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A_{ij}} \sum_{mn} A_{mn} B_{mn} = B_{ij}$$
$$\Rightarrow \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = B$$
$$\frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} F_1(A)^{\mathsf{T}} B F_2(A) C \end{bmatrix} \text{ with } F_1 \text{ and } F_2 \text{ both identity maps}$$
$$= \frac{\partial}{\partial F_1} \operatorname{Tr} \begin{bmatrix} F_1^{\mathsf{T}} B F_2 C \end{bmatrix} \frac{\partial F_1}{\partial A} + \frac{\partial}{\partial F_2} \operatorname{Tr} \begin{bmatrix} F_1^{\mathsf{T}} B F_2 C \end{bmatrix} \frac{\partial F_2}{\partial A}$$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \begin{bmatrix} \mathbf{x}^{\mathsf{T}} A \mathbf{y} \end{bmatrix} \text{ (scalars equal their own transpose and trace)}$$
$$\operatorname{Tr} [A] = \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} \end{bmatrix} \qquad \operatorname{Tr} [ABC] = \operatorname{Tr} [CAB] = \operatorname{Tr} [BCA]$$
$$\frac{\partial}{\partial A_{ij}} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A_{ij}} \sum_{mn} A_{mn} B_{mn} = B_{ij}$$
$$\Rightarrow \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = B$$
$$\frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B A \end{bmatrix} = \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} F_1(A)^{\mathsf{T}} B F_2(A) C \end{bmatrix} \text{ with } F_1 \text{ and } F_2 \text{ both identity maps}$$
$$= \frac{\partial}{\partial F_1} \operatorname{Tr} \begin{bmatrix} F_1^{\mathsf{T}} B F_2 C \end{bmatrix} \frac{\partial F_1}{\partial A} + \frac{\partial}{\partial F_2} \operatorname{Tr} \begin{bmatrix} CF_1^{\mathsf{T}} B F_2 \end{bmatrix} \frac{\partial F_2}{\partial A}$$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \begin{bmatrix} \mathbf{x}^{\mathsf{T}} A \mathbf{y} \end{bmatrix} \text{ (scalars equal their own transpose and trace)}$$
$$\operatorname{Tr} [A] = \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} \end{bmatrix} \qquad \operatorname{Tr} [ABC] = \operatorname{Tr} [CAB] = \operatorname{Tr} [BCA]$$
$$\frac{\partial}{\partial A_{ij}} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A_{ij}} \sum_{mn} A_{mn} B_{mn} = B_{ij}$$
$$\Rightarrow \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = B$$
$$\frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B A \end{bmatrix} = \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} F_1(A)^{\mathsf{T}} B F_2(A) C \end{bmatrix} \text{ with } F_1 \text{ and } F_2 \text{ both identity maps}$$
$$= \frac{\partial}{\partial F_1} \operatorname{Tr} \begin{bmatrix} F_1^{\mathsf{T}} B F_2 C \end{bmatrix} \frac{\partial F_1}{\partial A} + \frac{\partial}{\partial F_2} \operatorname{Tr} \begin{bmatrix} F_2^{\mathsf{T}} B^{\mathsf{T}} F_1 C^{\mathsf{T}} \end{bmatrix} \frac{\partial F_2}{\partial A}$$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \begin{bmatrix} \mathbf{x}^{\mathsf{T}} A \mathbf{y} \end{bmatrix} \text{ (scalars equal their own transpose and trace)}$$

$$\operatorname{Tr} [A] = \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} \end{bmatrix} \qquad \operatorname{Tr} [ABC] = \operatorname{Tr} [CAB] = \operatorname{Tr} [BCA]$$

$$\frac{\partial}{\partial A_{ij}} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A_{ij}} \sum_{mn} A_{mn} B_{mn} = B_{ij}$$

$$\Rightarrow \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = B$$

$$\frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B A \end{bmatrix} = \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} F_1(A)^{\mathsf{T}} B F_2(A) C \end{bmatrix} \text{ with } F_1 \text{ and } F_2 \text{ both identity maps}$$

$$= \frac{\partial}{\partial F_1} \operatorname{Tr} \begin{bmatrix} F_1^{\mathsf{T}} B F_2 C \end{bmatrix} \frac{\partial F_1}{\partial A} + \frac{\partial}{\partial F_2} \operatorname{Tr} \begin{bmatrix} F_2^{\mathsf{T}} B^{\mathsf{T}} F_1 C^{\mathsf{T}} \end{bmatrix} \frac{\partial F_2}{\partial A}$$

$$= BF_2 C + B^{\mathsf{T}} F_1 C^{\mathsf{T}}$$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \begin{bmatrix} \mathbf{x}^{\mathsf{T}} A \mathbf{y} \end{bmatrix} \text{ (scalars equal their own transpose and trace)}$$

$$\operatorname{Tr} [A] = \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} \end{bmatrix} \qquad \operatorname{Tr} [ABC] = \operatorname{Tr} [CAB] = \operatorname{Tr} [BCA]$$

$$\frac{\partial}{\partial A_{ij}} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A_{ij}} \sum_{mn} A_{mn} B_{mn} = B_{ij}$$

$$\Rightarrow \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = B$$

$$\frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B A \end{bmatrix} = \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} F_1(A)^{\mathsf{T}} B F_2(A) C \end{bmatrix} \text{ with } F_1 \text{ and } F_2 \text{ both identity maps}$$

$$= \frac{\partial}{\partial F_1} \operatorname{Tr} \begin{bmatrix} F_1^{\mathsf{T}} B F_2 C \end{bmatrix} \frac{\partial F_1}{\partial A} + \frac{\partial}{\partial F_2} \operatorname{Tr} \begin{bmatrix} F_2^{\mathsf{T}} B^{\mathsf{T}} F_1 C^{\mathsf{T}} \end{bmatrix} \frac{\partial F_2}{\partial A}$$

$$= BF_2 C + B^{\mathsf{T}} F_1 C^{\mathsf{T}} = BAC + B^{\mathsf{T}} A C^{\mathsf{T}}$$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \begin{bmatrix} \mathbf{x}^{\mathsf{T}} A \mathbf{y} \end{bmatrix} \text{ (scalars equal their own transpose and trace)}$$
$$\operatorname{Tr} [A] = \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} \end{bmatrix} \qquad \operatorname{Tr} [ABC] = \operatorname{Tr} [CAB] = \operatorname{Tr} [BCA]$$
$$\frac{\partial}{\partial A_{ij}} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A_{ij}} \sum_{mn} A_{mn} B_{mn} = B_{ij}$$
$$\Rightarrow \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = B$$
$$\frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = B$$
$$= \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} F_1(A)^{\mathsf{T}} B F_2(A) C \end{bmatrix} \text{ with } F_1 \text{ and } F_2 \text{ both identity maps}$$
$$= \frac{\partial}{\partial F_1} \operatorname{Tr} \begin{bmatrix} F_1^{\mathsf{T}} B F_2 C \end{bmatrix} \frac{\partial F_1}{\partial A} + \frac{\partial}{\partial F_2} \operatorname{Tr} \begin{bmatrix} F_2^{\mathsf{T}} B^{\mathsf{T}} F_1 C^{\mathsf{T}} \end{bmatrix} \frac{\partial F_2}{\partial A}$$
$$= BF_2 C + B^{\mathsf{T}} F_1 C^{\mathsf{T}} = BAC + B^{\mathsf{T}} A C^{\mathsf{T}}$$
$$\frac{\partial}{\partial A_{ij}} \log |A|$$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \begin{bmatrix} \mathbf{x}^{\mathsf{T}} A \mathbf{y} \end{bmatrix} \text{ (scalars equal their own transpose and trace)}$$
$$\operatorname{Tr} [A] = \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} \end{bmatrix} \qquad \operatorname{Tr} [ABC] = \operatorname{Tr} [CAB] = \operatorname{Tr} [BCA]$$
$$\frac{\partial}{\partial A_{ij}} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A_{ij}} \sum_{mn} A_{mn} B_{mn} = B_{ij}$$
$$\Rightarrow \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = B$$
$$\frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} F_1(A)^{\mathsf{T}} B F_2(A) C \end{bmatrix} \text{ with } F_1 \text{ and } F_2 \text{ both identity maps}$$
$$= \frac{\partial}{\partial F_1} \operatorname{Tr} \begin{bmatrix} F_1^{\mathsf{T}} B F_2 C \end{bmatrix} \frac{\partial F_1}{\partial A} + \frac{\partial}{\partial F_2} \operatorname{Tr} \begin{bmatrix} F_2^{\mathsf{T}} B^{\mathsf{T}} F_1 C^{\mathsf{T}} \end{bmatrix} \frac{\partial F_2}{\partial A}$$
$$= BF_2 C + B^{\mathsf{T}} F_1 C^{\mathsf{T}} = BAC + B^{\mathsf{T}} A C^{\mathsf{T}}$$
$$\frac{\partial}{\partial A_{ij}} \log |A| = \frac{1}{|A|} \frac{\partial}{\partial A_{ij}} |A|$$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \begin{bmatrix} \mathbf{x}^{\mathsf{T}} A \mathbf{y} \end{bmatrix} \text{ (scalars equal their own transpose and trace)}$$

$$\operatorname{Tr} [A] = \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} \end{bmatrix} \qquad \operatorname{Tr} [ABC] = \operatorname{Tr} [CAB] = \operatorname{Tr} [BCA]$$

$$\frac{\partial}{\partial A_{ij}} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A_{ij}} \sum_{mn} A_{mn} B_{mn} = B_{ij}$$

$$\Rightarrow \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = B$$

$$\frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B A C \end{bmatrix} = \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} F_1(A)^{\mathsf{T}} B F_2(A) C \end{bmatrix} \text{ with } F_1 \text{ and } F_2 \text{ both identity maps}$$

$$= \frac{\partial}{\partial F_1} \operatorname{Tr} \begin{bmatrix} F_1^{\mathsf{T}} B F_2 C \end{bmatrix} \frac{\partial F_1}{\partial A} + \frac{\partial}{\partial F_2} \operatorname{Tr} \begin{bmatrix} F_2^{\mathsf{T}} B^{\mathsf{T}} F_1 C^{\mathsf{T}} \end{bmatrix} \frac{\partial F_2}{\partial A}$$

$$= BF_2 C + B^{\mathsf{T}} F_1 C^{\mathsf{T}} = BAC + B^{\mathsf{T}} A C^{\mathsf{T}}$$

$$\frac{\partial}{\partial A_{ij}} \log |A| = \frac{1}{|A|} \frac{\partial}{\partial A_{ij}} \sum_{k} (-1)^{i+k} A_{ik} |[A]_{ik}|$$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \begin{bmatrix} \mathbf{x}^{\mathsf{T}} A \mathbf{y} \end{bmatrix} \text{ (scalars equal their own transpose and trace)}$$

$$\operatorname{Tr} [A] = \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} \end{bmatrix} \qquad \operatorname{Tr} [ABC] = \operatorname{Tr} [CAB] = \operatorname{Tr} [BCA]$$

$$\frac{\partial}{\partial A_{ij}} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A_{ij}} \sum_{mn} A_{mn} B_{mn} = B_{ij}$$

$$\Rightarrow \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = B$$

$$\frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B A C \end{bmatrix} = \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} F_1(A)^{\mathsf{T}} B F_2(A) C \end{bmatrix} \text{ with } F_1 \text{ and } F_2 \text{ both identity maps}$$

$$= \frac{\partial}{\partial F_1} \operatorname{Tr} \begin{bmatrix} F_1^{\mathsf{T}} B F_2 C \end{bmatrix} \frac{\partial F_1}{\partial A} + \frac{\partial}{\partial F_2} \operatorname{Tr} \begin{bmatrix} F_2^{\mathsf{T}} B^{\mathsf{T}} F_1 C^{\mathsf{T}} \end{bmatrix} \frac{\partial F_2}{\partial A}$$

$$= BF_2 C + B^{\mathsf{T}} F_1 C^{\mathsf{T}} = BAC + B^{\mathsf{T}} A C^{\mathsf{T}}$$

$$\frac{\partial}{\partial A_{ij}} \log |A| = \frac{1}{|A|} \frac{\partial}{\partial A_{ij}} \sum_{k} (-1)^{i+k} A_{ik} |[A]_{ik}| = \frac{1}{|A|} (-1)^{i+j} |[A]_{ij}|$$

$$\mathbf{x}^{\mathsf{T}} A \mathbf{y} = \mathbf{y}^{\mathsf{T}} A^{\mathsf{T}} \mathbf{x} = \operatorname{Tr} \begin{bmatrix} \mathbf{x}^{\mathsf{T}} A \mathbf{y} \end{bmatrix} \text{ (scalars equal their own transpose and trace)}$$

$$\operatorname{Tr} [A] = \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} \end{bmatrix} \quad \operatorname{Tr} [ABC] = \operatorname{Tr} [CAB] = \operatorname{Tr} [BCA]$$

$$\frac{\partial}{\partial A_{ij}} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A_{ij}} \sum_{mn} A_{mn} B_{mn} = B_{ij}$$

$$\Rightarrow \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = B$$

$$\frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} A^{\mathsf{T}} B \end{bmatrix} = \frac{\partial}{\partial A} \operatorname{Tr} \begin{bmatrix} F_1(A)^{\mathsf{T}} B F_2(A) C \end{bmatrix} \text{ with } F_1 \text{ and } F_2 \text{ both identity maps}$$

$$= \frac{\partial}{\partial F_1} \operatorname{Tr} \begin{bmatrix} F_1^{\mathsf{T}} B F_2 C \end{bmatrix} \frac{\partial F_1}{\partial A} + \frac{\partial}{\partial F_2} \operatorname{Tr} \begin{bmatrix} F_2^{\mathsf{T}} B^{\mathsf{T}} F_1 C^{\mathsf{T}} \end{bmatrix} \frac{\partial F_2}{\partial A}$$

$$= BF_2 C + B^{\mathsf{T}} F_1 C^{\mathsf{T}} = BAC + B^{\mathsf{T}} AC^{\mathsf{T}}$$

$$\frac{\partial}{\partial A_{ij}} \log |A| = \frac{1}{|A|} \frac{\partial}{\partial A_{ij}} \sum_{k} (-1)^{i+k} A_{ik} |[A]_{ik}| = \frac{1}{|A|} (-1)^{i+j} |[A]_{ij}|$$

$$\Rightarrow \frac{\partial}{\partial A} \log |A| = (A^{-1})^{\mathsf{T}}$$

$$\frac{\partial(-\ell)}{\partial \mu} = \frac{\partial}{\partial \mu} \left[\frac{N}{2} \log |2\pi\Sigma| + \frac{1}{2} \sum_{n} (\mathbf{x}_{n} - \mu)^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \mu) \right]$$

$$\begin{split} \frac{\partial(-\ell)}{\partial\mu} &= \frac{\partial}{\partial\mu} \left[\frac{N}{2} \log |2\pi\Sigma| + \frac{1}{2} \sum_{n} (\mathbf{x}_{n} - \mu)^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \mu) \right] \\ &= \frac{1}{2} \sum_{n} \frac{\partial}{\partial\mu} \left[(\mathbf{x}_{n} - \mu)^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \mu) \right] \end{split}$$

$$\begin{aligned} \frac{\partial(-\ell)}{\partial\mu} &= \frac{\partial}{\partial\mu} \left[\frac{N}{2} \log |2\pi\Sigma| + \frac{1}{2} \sum_{n} (\mathbf{x}_{n} - \mu)^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \mu) \right] \\ &= \frac{1}{2} \sum_{n} \frac{\partial}{\partial\mu} \left[(\mathbf{x}_{n} - \mu)^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \mu) \right] \\ &= \frac{1}{2} \sum_{n} \frac{\partial}{\partial\mu} \left[\mathbf{x}_{n}^{\mathsf{T}} \Sigma^{-1} \mathbf{x}_{n} + \mu^{\mathsf{T}} \Sigma^{-1} \mu - 2\mu^{\mathsf{T}} \Sigma^{-1} \mathbf{x}_{n} \right] \end{aligned}$$

$$\frac{\partial(-\ell)}{\partial\mu} = \frac{\partial}{\partial\mu} \left[\frac{N}{2} \log |2\pi\Sigma| + \frac{1}{2} \sum_{n} (\mathbf{x}_{n} - \mu)^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \mu) \right]$$
$$= \frac{1}{2} \sum_{n} \frac{\partial}{\partial\mu} \left[(\mathbf{x}_{n} - \mu)^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \mu) \right]$$
$$= \frac{1}{2} \sum_{n} \frac{\partial}{\partial\mu} \left[\mathbf{x}_{n}^{\mathsf{T}} \Sigma^{-1} \mathbf{x}_{n} + \mu^{\mathsf{T}} \Sigma^{-1} \mu - 2\mu^{\mathsf{T}} \Sigma^{-1} \mathbf{x}_{n} \right]$$
$$= \frac{1}{2} \sum_{n} \frac{\partial}{\partial\mu} \left[\mu^{\mathsf{T}} \Sigma^{-1} \mu \right] - 2\frac{\partial}{\partial\mu} \left[\mu^{\mathsf{T}} \Sigma^{-1} \mathbf{x}_{n} \right]$$

$$\begin{aligned} \frac{\partial(-\ell)}{\partial\mu} &= \frac{\partial}{\partial\mu} \left[\frac{N}{2} \log |2\pi\Sigma| + \frac{1}{2} \sum_{n} (\mathbf{x}_{n} - \mu)^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \mu) \right] \\ &= \frac{1}{2} \sum_{n} \frac{\partial}{\partial\mu} \left[(\mathbf{x}_{n} - \mu)^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \mu) \right] \\ &= \frac{1}{2} \sum_{n} \frac{\partial}{\partial\mu} \left[\mathbf{x}_{n}^{\mathsf{T}} \Sigma^{-1} \mathbf{x}_{n} + \mu^{\mathsf{T}} \Sigma^{-1} \mu - 2\mu^{\mathsf{T}} \Sigma^{-1} \mathbf{x}_{n} \right] \\ &= \frac{1}{2} \sum_{n} \frac{\partial}{\partial\mu} \left[\mu^{\mathsf{T}} \Sigma^{-1} \mu \right] - 2\frac{\partial}{\partial\mu} \left[\mu^{\mathsf{T}} \Sigma^{-1} \mathbf{x}_{n} \right] \\ &= \frac{1}{2} \sum_{n} \left[2\Sigma^{-1} \mu - 2\Sigma^{-1} \mathbf{x}_{n} \right] \end{aligned}$$

$$\frac{\partial(-\ell)}{\partial\mu} = \frac{\partial}{\partial\mu} \left[\frac{N}{2} \log |2\pi\Sigma| + \frac{1}{2} \sum_{n} (\mathbf{x}_{n} - \mu)^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \mu) \right]$$
$$= \frac{1}{2} \sum_{n} \frac{\partial}{\partial\mu} \left[(\mathbf{x}_{n} - \mu)^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \mu) \right]$$
$$= \frac{1}{2} \sum_{n} \frac{\partial}{\partial\mu} \left[\mathbf{x}_{n}^{\mathsf{T}} \Sigma^{-1} \mathbf{x}_{n} + \mu^{\mathsf{T}} \Sigma^{-1} \mu - 2\mu^{\mathsf{T}} \Sigma^{-1} \mathbf{x}_{n} \right]$$
$$= \frac{1}{2} \sum_{n} \frac{\partial}{\partial\mu} \left[\mu^{\mathsf{T}} \Sigma^{-1} \mu \right] - 2 \frac{\partial}{\partial\mu} \left[\mu^{\mathsf{T}} \Sigma^{-1} \mathbf{x}_{n} \right]$$
$$= \frac{1}{2} \sum_{n} \left[2\Sigma^{-1} \mu - 2\Sigma^{-1} \mathbf{x}_{n} \right]$$
$$= N\Sigma^{-1} \mu - \Sigma^{-1} \sum_{n} \mathbf{x}_{n}$$

$$\frac{\partial(-\ell)}{\partial\mu} = \frac{\partial}{\partial\mu} \left[\frac{N}{2} \log |2\pi\Sigma| + \frac{1}{2} \sum_{n} (\mathbf{x}_{n} - \mu)^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \mu) \right]$$
$$= \frac{1}{2} \sum_{n} \frac{\partial}{\partial\mu} \left[(\mathbf{x}_{n} - \mu)^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \mu) \right]$$
$$= \frac{1}{2} \sum_{n} \frac{\partial}{\partial\mu} \left[\mathbf{x}_{n}^{\mathsf{T}} \Sigma^{-1} \mathbf{x}_{n} + \mu^{\mathsf{T}} \Sigma^{-1} \mu - 2\mu^{\mathsf{T}} \Sigma^{-1} \mathbf{x}_{n} \right]$$
$$= \frac{1}{2} \sum_{n} \frac{\partial}{\partial\mu} \left[\mu^{\mathsf{T}} \Sigma^{-1} \mu \right] - 2 \frac{\partial}{\partial\mu} \left[\mu^{\mathsf{T}} \Sigma^{-1} \mathbf{x}_{n} \right]$$
$$= \frac{1}{2} \sum_{n} \left[2\Sigma^{-1} \mu - 2\Sigma^{-1} \mathbf{x}_{n} \right]$$
$$= N\Sigma^{-1} \mu - \Sigma^{-1} \sum_{n} \mathbf{x}_{n}$$

 $= 0 \Rightarrow \hat{\mu} = \frac{1}{N} \sum_{n} \mathbf{x}_{n}$



$$\frac{\partial(-\ell)}{\partial \boldsymbol{\Sigma}^{-1}} = \frac{\partial}{\partial \boldsymbol{\Sigma}^{-1}} \left[\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}| + \frac{1}{2} \sum_{n} (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x}_{n} - \boldsymbol{\mu}) \right]$$

$$\frac{\partial(-\ell)}{\partial \Sigma^{-1}} = \frac{\partial}{\partial \Sigma^{-1}} \left[\frac{N}{2} \log |2\pi\Sigma| + \frac{1}{2} \sum_{n} (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \boldsymbol{\mu}) \right]$$
$$= \frac{\partial}{\partial \Sigma^{-1}} \left[\frac{N}{2} \log |2\piI| \right] - \frac{\partial}{\partial \Sigma^{-1}} \left[\frac{N}{2} \log |\Sigma^{-1}| \right]$$
$$+ \frac{1}{2} \sum_{n} \frac{\partial}{\partial \Sigma^{-1}} \left[(\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \boldsymbol{\mu}) \right]$$

$$\frac{\partial(-\ell)}{\partial \Sigma^{-1}} = \frac{\partial}{\partial \Sigma^{-1}} \left[\frac{N}{2} \log |2\pi\Sigma| + \frac{1}{2} \sum_{n} (\mathbf{x}_{n} - \mu)^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \mu) \right]$$
$$= \frac{\partial}{\partial \Sigma^{-1}} \left[\frac{N}{2} \log |2\piI| \right] - \frac{\partial}{\partial \Sigma^{-1}} \left[\frac{N}{2} \log |\Sigma^{-1}| \right]$$
$$+ \frac{1}{2} \sum_{n} \frac{\partial}{\partial \Sigma^{-1}} \left[(\mathbf{x}_{n} - \mu)^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \mu) \right]$$
$$= -\frac{N}{2} \Sigma^{\mathsf{T}} + \frac{1}{2} \sum_{n} (\mathbf{x}_{n} - \mu) (\mathbf{x}_{n} - \mu)^{\mathsf{T}}$$

$$\frac{\partial(-\ell)}{\partial \Sigma^{-1}} = \frac{\partial}{\partial \Sigma^{-1}} \left[\frac{N}{2} \log |2\pi\Sigma| + \frac{1}{2} \sum_{n} (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \boldsymbol{\mu}) \right]$$
$$= \frac{\partial}{\partial \Sigma^{-1}} \left[\frac{N}{2} \log |2\piI| \right] - \frac{\partial}{\partial \Sigma^{-1}} \left[\frac{N}{2} \log |\Sigma^{-1}| \right]$$
$$+ \frac{1}{2} \sum_{n} \frac{\partial}{\partial \Sigma^{-1}} \left[(\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathsf{T}} \Sigma^{-1} (\mathbf{x}_{n} - \boldsymbol{\mu}) \right]$$
$$= -\frac{N}{2} \Sigma^{\mathsf{T}} + \frac{1}{2} \sum_{n} (\mathbf{x}_{n} - \boldsymbol{\mu}) (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathsf{T}}$$
$$= 0 \Rightarrow \hat{\Sigma} = \frac{1}{N} \sum_{n} (\mathbf{x}_{n} - \boldsymbol{\mu}) (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathsf{T}}$$

Equivalences



Multivariate Linear Regression

The relationship between variables can also be modelled as a conditional distribution.



- data $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{y}_1) \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$
- each \mathbf{x}_i (\mathbf{y}_i) is a vector of D_x (D_y) features,
- **y**_{*i*} is conditionally independent of all else, given **x**_{*i*}.

Multivariate Linear Regression

The relationship between variables can also be modelled as a conditional distribution.



- data $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{y}_1) \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$
- each \mathbf{x}_i (\mathbf{y}_i) is a vector of D_x (D_y) features,
- **y**_{*i*} is conditionally independent of all else, given **x**_{*i*}.

A simple form of supervised (predictive) learning: model ${\bf y}$ as a **linear** function of ${\bf x}$, with **Gaussian** noise.
Multivariate Linear Regression

The relationship between variables can also be modelled as a conditional distribution.



- data $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{y}_1) \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$
- each \mathbf{x}_i (\mathbf{y}_i) is a vector of D_x (D_y) features,
- **y**_{*i*} is conditionally independent of all else, given **x**_{*i*}.

A simple form of supervised (predictive) learning: model **y** as a **linear** function of **x**, with **Gaussian** noise.

$$p(\mathbf{y}|\mathbf{x}, \mathsf{W}, \boldsymbol{\Sigma}_{y}) = |2\pi\boldsymbol{\Sigma}_{y}|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(\mathbf{y} - \mathsf{W}\mathbf{x})^{\mathsf{T}}\boldsymbol{\Sigma}_{y}^{-1}(\mathbf{y} - \mathsf{W}\mathbf{x})\right\}$$

$$\begin{split} \ell &= \sum_{i} \log p(\mathbf{y}_{i} | \mathbf{x}_{i}, \mathsf{W}, \boldsymbol{\Sigma}_{y}) \\ &= -\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}_{y}| - \frac{1}{2} \sum_{i} (\mathbf{y}_{i} - \mathsf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathsf{W}\mathbf{x}_{i}) \end{split}$$

$$\begin{split} \ell &= \sum_{i} \log p(\mathbf{y}_{i} | \mathbf{x}_{i}, \mathbf{W}, \boldsymbol{\Sigma}_{y}) \\ &= -\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}_{y}| - \frac{1}{2} \sum_{i} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i}) \\ \frac{\partial (-\ell)}{\partial \mathbf{W}} \end{split}$$

$$\begin{split} \ell &= \sum_{i} \log p(\mathbf{y}_{i} | \mathbf{x}_{i}, \mathbf{W}, \boldsymbol{\Sigma}_{y}) \\ &= -\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}_{y}| - \frac{1}{2} \sum_{i} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i}) \\ \frac{\partial (-\ell)}{\partial \mathbf{W}} &= \frac{\partial}{\partial \mathbf{W}} \left[\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}_{y}| + \frac{1}{2} \sum_{i} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i}) \right] \end{split}$$

$$\begin{split} \ell &= \sum_{i} \log p(\mathbf{y}_{i} | \mathbf{x}_{i}, \mathbf{W}, \boldsymbol{\Sigma}_{y}) \\ &= -\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}_{y}| - \frac{1}{2} \sum_{i} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i}) \\ \frac{\partial (-\ell)}{\partial \mathbf{W}} &= \frac{\partial}{\partial \mathbf{W}} \left[\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}_{y}| + \frac{1}{2} \sum_{i} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i}) \right] \\ &= \frac{1}{2} \sum_{i} \frac{\partial}{\partial \mathbf{W}} \left[(\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i}) \right] \end{split}$$

$$\ell = \sum_{i} \log p(\mathbf{y}_{i} | \mathbf{x}_{i}, \mathbf{W}, \boldsymbol{\Sigma}_{y})$$

$$= -\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}_{y}| - \frac{1}{2} \sum_{i} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})$$

$$\frac{\partial(-\ell)}{\partial \mathbf{W}} = \frac{\partial}{\partial \mathbf{W}} \left[\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}_{y}| + \frac{1}{2} \sum_{i} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i}) \right]$$

$$= \frac{1}{2} \sum_{i} \frac{\partial}{\partial \mathbf{W}} \left[(\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i}) \right]$$

$$= \frac{1}{2} \sum_{i} \frac{\partial}{\partial \mathbf{W}} \left[\mathbf{y}_{i}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathbf{y}_{i} + \mathbf{x}_{i}^{\mathsf{T}} \mathbf{W}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathbf{W}\mathbf{x}_{i} - 2\mathbf{x}_{i}^{\mathsf{T}} \mathbf{W}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathbf{y}_{i} \right]$$

$$\ell = \sum_{i} \log p(\mathbf{y}_{i} | \mathbf{x}_{i}, \mathbf{W}, \boldsymbol{\Sigma}_{y})$$

$$= -\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}_{y}| - \frac{1}{2} \sum_{i} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})$$

$$\frac{\partial(-\ell)}{\partial \mathbf{W}} = \frac{\partial}{\partial \mathbf{W}} \left[\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}_{y}| + \frac{1}{2} \sum_{i} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i}) \right]$$

$$= \frac{1}{2} \sum_{i} \frac{\partial}{\partial \mathbf{W}} \left[(\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i}) \right]$$

$$= \frac{1}{2} \sum_{i} \frac{\partial}{\partial \mathbf{W}} \left[\mathbf{y}_{i}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathbf{y}_{i} + \mathbf{x}_{i}^{\mathsf{T}} \mathbf{W}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathbf{W}\mathbf{x}_{i} - 2\mathbf{x}_{i}^{\mathsf{T}} \mathbf{W}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathbf{y}_{i} \right]$$

$$= \frac{1}{2} \sum_{i} \left[\frac{\partial}{\partial \mathbf{W}} \operatorname{Tr} \left[\mathbf{W}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathbf{W} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathsf{T}} \right] - 2 \frac{\partial}{\partial \mathbf{W}} \operatorname{Tr} \left[\mathbf{W}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathbf{y}_{i} \mathbf{x}_{i}^{\mathsf{T}} \right]$$

$$\begin{split} \ell &= \sum_{i} \log p(\mathbf{y}_{i} | \mathbf{x}_{i}, \mathbf{W}, \boldsymbol{\Sigma}_{y}) \\ &= -\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}_{y}| - \frac{1}{2} \sum_{i} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i}) \\ \frac{\partial(-\ell)}{\partial \mathsf{W}} &= \frac{\partial}{\partial \mathsf{W}} \left[\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}_{y}| + \frac{1}{2} \sum_{i} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i}) \right] \\ &= \frac{1}{2} \sum_{i} \frac{\partial}{\partial \mathsf{W}} \left[(\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i}) \right] \\ &= \frac{1}{2} \sum_{i} \frac{\partial}{\partial \mathsf{W}} \left[\mathbf{y}_{i}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathbf{y}_{i} + \mathbf{x}_{i}^{\mathsf{T}} \mathsf{W}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathsf{W} \mathbf{x}_{i} - 2\mathbf{x}_{i}^{\mathsf{T}} \mathsf{W}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathbf{y}_{i} \right] \\ &= \frac{1}{2} \sum_{i} \left[\frac{\partial}{\partial \mathsf{W}} \operatorname{Tr} \left[\mathsf{W}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathsf{W} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathsf{T}} \right] - 2 \frac{\partial}{\partial \mathsf{W}} \operatorname{Tr} \left[\mathsf{W}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathbf{y}_{i} \mathbf{x}_{i}^{\mathsf{T}} \right] \\ &= \frac{1}{2} \sum_{i} \left[2 \boldsymbol{\Sigma}_{y}^{-1} \mathsf{W} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathsf{T}} - 2 \boldsymbol{\Sigma}_{y}^{-1} \mathbf{y}_{i} \mathbf{x}_{i}^{\mathsf{T}} \right] \end{split}$$

$$\ell = \sum_{i} \log p(\mathbf{y}_{i} | \mathbf{x}_{i}, \mathbf{W}, \boldsymbol{\Sigma}_{y})$$

$$= -\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}_{y}| - \frac{1}{2} \sum_{i} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})$$

$$\frac{\partial(-\ell)}{\partial \mathbf{W}} = \frac{\partial}{\partial \mathbf{W}} \left[\frac{N}{2} \log |2\pi\boldsymbol{\Sigma}_{y}| + \frac{1}{2} \sum_{i} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i}) \right]$$

$$= \frac{1}{2} \sum_{i} \frac{\partial}{\partial \mathbf{W}} \left[(\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i})^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} (\mathbf{y}_{i} - \mathbf{W}\mathbf{x}_{i}) \right]$$

$$= \frac{1}{2} \sum_{i} \frac{\partial}{\partial \mathbf{W}} \left[\mathbf{y}_{i}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathbf{y}_{i} + \mathbf{x}_{i}^{\mathsf{T}} \mathbf{W}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathbf{W}\mathbf{x}_{i} - 2\mathbf{x}_{i}^{\mathsf{T}} \mathbf{W}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathbf{y}_{i} \right]$$

$$= \frac{1}{2} \sum_{i} \left[\frac{\partial}{\partial \mathbf{W}} \operatorname{Tr} \left[\mathbf{W}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathbf{W} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathsf{T}} \right] - 2 \frac{\partial}{\partial \mathbf{W}} \operatorname{Tr} \left[\mathbf{W}^{\mathsf{T}} \boldsymbol{\Sigma}_{y}^{-1} \mathbf{y}_{i} \mathbf{x}_{i}^{\mathsf{T}} \right]$$

$$= \frac{1}{2} \sum_{i} \left[2\boldsymbol{\Sigma}_{y}^{-1} \mathbf{W} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathsf{T}} - 2\boldsymbol{\Sigma}_{y}^{-1} \mathbf{y}_{i} \mathbf{x}_{i}^{\mathsf{T}} \right]$$

$$= 0 \Rightarrow \widehat{\mathbf{W}} = \sum_{i} \mathbf{y}_{i} \mathbf{x}_{i}^{\mathsf{T}} \left(\sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathsf{T}} \right)^{-1}$$

Let y_i be scalar (so that W is a row vector) and write **w** for the column vector of weights.

Let y_i be scalar (so that W is a row vector) and write **w** for the column vector of weights. A conjugate prior for **w** is

$$P(\mathbf{w}|\mathbf{A}) = \mathcal{N}\left(\mathbf{0}, \mathbf{A}^{-1}\right)$$

Let y_i be scalar (so that W is a row vector) and write **w** for the column vector of weights. A conjugate prior for **w** is

$$P(\mathbf{w}|\mathsf{A}) = \mathcal{N}\left(\mathbf{0}, \mathsf{A}^{-1}\right)$$

$$\log P(\mathbf{w}|\mathcal{D}, \mathsf{A}, \sigma_y) = \log P(\mathcal{D}|\mathbf{w}, \mathsf{A}, \sigma_y) + \log P(\mathbf{w}|\mathsf{A}, \sigma_y) - \log P(\mathcal{D}|\mathsf{A}, \sigma_y)$$

Let y_i be scalar (so that W is a row vector) and write **w** for the column vector of weights. A conjugate prior for **w** is

$$P(\mathbf{w}|\mathsf{A}) = \mathcal{N}\left(\mathbf{0}, \mathsf{A}^{-1}\right)$$

$$\log P(\mathbf{w}|\mathcal{D}, \mathsf{A}, \sigma_y) = \log P(\mathcal{D}|\mathbf{w}, \mathsf{A}, \sigma_y) + \log P(\mathbf{w}|\mathsf{A}, \sigma_y) - \log P(\mathcal{D}|\mathsf{A}, \sigma_y)$$
$$= -\frac{1}{2}\mathbf{w}^{\mathsf{T}}\mathsf{A}\mathbf{w} - \frac{1}{2}\sum_{i}(y_i - \mathbf{w}^{\mathsf{T}}\mathbf{x}_i)^2\sigma_y^{-2} + \text{const}$$

Let y_i be scalar (so that W is a row vector) and write **w** for the column vector of weights. A conjugate prior for **w** is

$$P(\mathbf{w}|\mathbf{A}) = \mathcal{N}\left(\mathbf{0}, \mathbf{A}^{-1}\right)$$

$$\log P(\mathbf{w}|\mathcal{D}, \mathbf{A}, \sigma_y) = \log P(\mathcal{D}|\mathbf{w}, \mathbf{A}, \sigma_y) + \log P(\mathbf{w}|\mathbf{A}, \sigma_y) - \log P(\mathcal{D}|\mathbf{A}, \sigma_y)$$
$$= -\frac{1}{2}\mathbf{w}^{\mathsf{T}}\mathbf{A}\mathbf{w} - \frac{1}{2}\sum_{i}(y_i - \mathbf{w}^{\mathsf{T}}\mathbf{x}_i)^2 \sigma_y^{-2} + \text{const}$$
$$= -\frac{1}{2}\mathbf{w}^{\mathsf{T}}(\mathbf{A} + \sigma_y^{-2}\sum_{i}\mathbf{x}_i\mathbf{x}_i^{\mathsf{T}})\mathbf{w} + \mathbf{w}^{\mathsf{T}}\sum_{i}(y_i\mathbf{x}_i)\sigma_y^{-2} + \text{const}$$

Let y_i be scalar (so that W is a row vector) and write **w** for the column vector of weights. A conjugate prior for **w** is

$$P(\mathbf{w}|\mathbf{A}) = \mathcal{N}\left(\mathbf{0}, \mathbf{A}^{-1}\right)$$

$$\log P(\mathbf{w}|\mathcal{D}, \mathbf{A}, \sigma_{y}) = \log P(\mathcal{D}|\mathbf{w}, \mathbf{A}, \sigma_{y}) + \log P(\mathbf{w}|\mathbf{A}, \sigma_{y}) - \log P(\mathcal{D}|\mathbf{A}, \sigma_{y})$$
$$= -\frac{1}{2} \mathbf{w}^{\mathsf{T}} \mathbf{A} \mathbf{w} - \frac{1}{2} \sum_{i} (y_{i} - \mathbf{w}^{\mathsf{T}} \mathbf{x}_{i})^{2} \sigma_{y}^{-2} + \text{const}$$
$$= -\frac{1}{2} \mathbf{w}^{\mathsf{T}} \underbrace{(\mathbf{A} + \sigma_{y}^{-2} \sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathsf{T}})}_{\boldsymbol{\Sigma}_{w}^{-1}} \mathbf{w} + \mathbf{w}^{\mathsf{T}} \sum_{i} (y_{i} \mathbf{x}_{i}) \sigma_{y}^{-2} + \text{const}$$

Let y_i be scalar (so that W is a row vector) and write **w** for the column vector of weights. A conjugate prior for **w** is

$$P(\mathbf{w}|\mathbf{A}) = \mathcal{N}\left(\mathbf{0}, \mathbf{A}^{-1}\right)$$

$$\log P(\mathbf{w}|\mathcal{D}, \mathbf{A}, \sigma_y) = \log P(\mathcal{D}|\mathbf{w}, \mathbf{A}, \sigma_y) + \log P(\mathbf{w}|\mathbf{A}, \sigma_y) - \log P(\mathcal{D}|\mathbf{A}, \sigma_y)$$
$$= -\frac{1}{2} \mathbf{w}^{\mathsf{T}} \mathbf{A} \mathbf{w} - \frac{1}{2} \sum_{i} (y_i - \mathbf{w}^{\mathsf{T}} \mathbf{x}_i)^2 \sigma_y^{-2} + \text{const}$$
$$= -\frac{1}{2} \mathbf{w}^{\mathsf{T}} \underbrace{(\mathbf{A} + \sigma_y^{-2} \sum_{i} \mathbf{x}_i \mathbf{x}_i^{\mathsf{T}})}_{\boldsymbol{\Sigma}_w^{-1}} \mathbf{w} + \mathbf{w}^{\mathsf{T}} \sum_{i} (y_i \mathbf{x}_i) \sigma_y^{-2} + \text{const}$$
$$= -\frac{1}{2} \mathbf{w}^{\mathsf{T}} \boldsymbol{\Sigma}_w^{-1} \mathbf{w} + \mathbf{w}^{\mathsf{T}} \boldsymbol{\Sigma}_w^{-1} \boldsymbol{\Sigma}_w \sum_{i} (y_i \mathbf{x}_i) \sigma_y^{-2} + \text{const}$$

Let y_i be scalar (so that W is a row vector) and write **w** for the column vector of weights. A conjugate prior for **w** is

$$P(\mathbf{w}|\mathbf{A}) = \mathcal{N}\left(\mathbf{0}, \mathbf{A}^{-1}\right)$$

$$\log P(\mathbf{w}|\mathcal{D}, \mathbf{A}, \sigma_y) = \log P(\mathcal{D}|\mathbf{w}, \mathbf{A}, \sigma_y) + \log P(\mathbf{w}|\mathbf{A}, \sigma_y) - \log P(\mathcal{D}|\mathbf{A}, \sigma_y)$$
$$= -\frac{1}{2} \mathbf{w}^{\mathsf{T}} \mathbf{A} \mathbf{w} - \frac{1}{2} \sum_{i} (y_i - \mathbf{w}^{\mathsf{T}} \mathbf{x}_i)^2 \sigma_y^{-2} + \text{const}$$
$$= -\frac{1}{2} \mathbf{w}^{\mathsf{T}} \underbrace{(\mathbf{A} + \sigma_y^{-2} \sum_{i} \mathbf{x}_i \mathbf{x}_i^{\mathsf{T}}) \mathbf{w}}_{\sum_{i=1}^{i-1}} + \mathbf{w}^{\mathsf{T}} \sum_{i} (y_i \mathbf{x}_i) \sigma_y^{-2} + \text{const}$$
$$= -\frac{1}{2} \mathbf{w}^{\mathsf{T}} \sum_{w}^{-1} \mathbf{w} + \mathbf{w}^{\mathsf{T}} \sum_{w}^{-1} \sum_{i} (y_i \mathbf{x}_i) \sigma_y^{-2} + \text{const}$$

Let y_i be scalar (so that W is a row vector) and write **w** for the column vector of weights. A conjugate prior for **w** is

$$P(\mathbf{w}|\mathbf{A}) = \mathcal{N}\left(\mathbf{0}, \mathbf{A}^{-1}\right)$$

$$\log P(\mathbf{w}|\mathcal{D}, \mathbf{A}, \sigma_y) = \log P(\mathcal{D}|\mathbf{w}, \mathbf{A}, \sigma_y) + \log P(\mathbf{w}|\mathbf{A}, \sigma_y) - \log P(\mathcal{D}|\mathbf{A}, \sigma_y)$$
$$= -\frac{1}{2} \mathbf{w}^{\mathsf{T}} \mathbf{A} \mathbf{w} - \frac{1}{2} \sum_{i} (y_i - \mathbf{w}^{\mathsf{T}} \mathbf{x}_i)^2 \sigma_y^{-2} + \text{const}$$
$$= -\frac{1}{2} \mathbf{w}^{\mathsf{T}} \underbrace{(\mathbf{A} + \sigma_y^{-2} \sum_{i} \mathbf{x}_i \mathbf{x}_i^{\mathsf{T}}) \mathbf{w}}_{\boldsymbol{\Sigma}_w^{-1}} + \mathbf{w}^{\mathsf{T}} \sum_{i} (y_i \mathbf{x}_i) \sigma_y^{-2} + \text{const}$$
$$= -\frac{1}{2} \mathbf{w}^{\mathsf{T}} \boldsymbol{\Sigma}_w^{-1} \mathbf{w} + \mathbf{w}^{\mathsf{T}} \boldsymbol{\Sigma}_w^{-1} \underbrace{\boldsymbol{\Sigma}_w \sum_{i} (y_i \mathbf{x}_i) \sigma_y^{-2}}_{\boldsymbol{\mu}_w} + \text{const}$$
$$= \log \mathcal{N} \left(\boldsymbol{\Sigma}_w \sum_{i} (y_i \mathbf{x}_i) \sigma_y^{-2}, \boldsymbol{\Sigma}_w \right)$$

As the posterior is Gaussian, the MAP and posterior mean weights are the same:

$$\mathbf{w}^{\text{MAP}} = \underbrace{\left(\mathbf{A} + \frac{\sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\text{T}}}{\sigma_{y}^{2}}\right)^{-1}}_{\Sigma_{w}} \underbrace{\frac{\sum_{i} y_{i} \mathbf{x}_{i}}{\sigma_{y}^{2}}}_{\Sigma_{w}}$$

As the posterior is Gaussian, the MAP and posterior mean weights are the same:

$$\mathbf{w}^{\text{MAP}} = \underbrace{\left(A + \frac{\sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\text{T}}}{\sigma_{y}^{2}}\right)^{-1}}_{\Sigma_{w}} \underbrace{\frac{\sum_{i} y_{i} \mathbf{x}_{i}}{\sigma_{y}^{2}}}_{\Sigma_{w}} = \left(A\sigma_{y}^{2} + \sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\text{T}}\right)^{-1} \sum_{i} y_{i} \mathbf{x}_{i}$$

As the posterior is Gaussian, the MAP and posterior mean weights are the same:

$$\mathbf{w}^{\text{MAP}} = \underbrace{\left(A + \frac{\sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\text{T}}}{\sigma_{y}^{2}}\right)^{-1}}_{\Sigma_{w}} \underbrace{\frac{\sum_{i} y_{i} \mathbf{x}_{i}}{\sigma_{y}^{2}}}_{\Sigma_{w}} = \left(A\sigma_{y}^{2} + \sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\text{T}}\right)^{-1} \sum_{i} y_{i} \mathbf{x}_{i}$$

$$\mathbf{w}^{\mathrm{ML}} = \widehat{\mathbf{W}}^{\mathrm{T}} = \left(\sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathrm{T}}\right)^{-1} \sum_{i} y_{i} \mathbf{x}_{i}$$

As the posterior is Gaussian, the MAP and posterior mean weights are the same:

$$\mathbf{w}^{\text{MAP}} = \underbrace{\left(A + \frac{\sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\text{T}}}{\sigma_{y}^{2}}\right)^{-1}}_{\Sigma_{w}} \underbrace{\frac{\sum_{i} y_{i} \mathbf{x}_{i}}{\sigma_{y}^{2}}}_{\Sigma_{w}} = \left(A\sigma_{y}^{2} + \sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\text{T}}\right)^{-1} \sum_{i} y_{i} \mathbf{x}_{i}$$

$$\mathbf{w}^{\mathrm{ML}} = \widehat{\mathbf{W}}^{\mathrm{T}} = \left(\sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathrm{T}}\right)^{-1} \sum_{i} y_{i} \mathbf{x}_{i}$$

As the posterior is Gaussian, the MAP and posterior mean weights are the same:

$$\mathbf{w}^{\text{MAP}} = \underbrace{\left(A + \frac{\sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\text{T}}}{\sigma_{y}^{2}}\right)^{-1}}_{\Sigma_{w}} \underbrace{\frac{\sum_{i} y_{i} \mathbf{x}_{i}}{\sigma_{y}^{2}}}_{\Sigma_{w}} = \left(A\sigma_{y}^{2} + \sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\text{T}}\right)^{-1} \sum_{i} y_{i} \mathbf{x}_{i}$$

Compare this to the (transposed) ML weight vector for scalar outputs:

$$\mathbf{w}^{\mathrm{ML}} = \widehat{\mathbf{W}}^{\mathrm{T}} = \left(\sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathrm{T}}\right)^{-1} \sum_{i} y_{i} \mathbf{x}_{i}$$

> The prior acts to "inflate" the apparent covariance of inputs.

As the posterior is Gaussian, the MAP and posterior mean weights are the same:

$$\mathbf{w}^{\text{MAP}} = \underbrace{\left(A + \frac{\sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\text{T}}}{\sigma_{y}^{2}}\right)^{-1}}_{\Sigma_{w}} \underbrace{\frac{\sum_{i} y_{i} \mathbf{x}_{i}}{\sigma_{y}^{2}}}_{\Sigma_{w}} = \left(A\sigma_{y}^{2} + \sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\text{T}}\right)^{-1} \sum_{i} y_{i} \mathbf{x}_{i}$$

$$\mathbf{w}^{\mathrm{ML}} = \widehat{\mathbf{W}}^{\mathrm{T}} = \left(\sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathrm{T}}\right)^{-1} \sum_{i} y_{i} \mathbf{x}_{i}$$

- The prior acts to "inflate" the apparent covariance of inputs.
- As A is positive (semi)definite, shrinks the weights towards the prior mean (here 0).

As the posterior is Gaussian, the MAP and posterior mean weights are the same:

$$\mathbf{w}^{\text{MAP}} = \underbrace{\left(A + \frac{\sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\text{T}}}{\sigma_{y}^{2}}\right)^{-1}}_{\Sigma_{w}} \underbrace{\frac{\sum_{i} y_{i} \mathbf{x}_{i}}{\sigma_{y}^{2}}}_{\Sigma_{w}} = \left(A\sigma_{y}^{2} + \sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\text{T}}\right)^{-1} \sum_{i} y_{i} \mathbf{x}_{i}$$

$$\mathbf{w}^{\mathrm{ML}} = \widehat{\mathbf{W}}^{\mathrm{T}} = \left(\sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathrm{T}}\right)^{-1} \sum_{i} y_{i} \mathbf{x}_{i}$$

- The prior acts to "inflate" the apparent covariance of inputs.
- As A is positive (semi)definite, shrinks the weights towards the prior mean (here 0).
- If $A = \alpha I$ this is known as the ridge regression estimator.

As the posterior is Gaussian, the MAP and posterior mean weights are the same:

$$\mathbf{w}^{\mathsf{MAP}} = \underbrace{\left(\mathbf{A} + \frac{\sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathsf{T}}}{\sigma_{y}^{2}}\right)^{-1}}_{\Sigma_{w}} \underbrace{\frac{\sum_{i} y_{i} \mathbf{x}_{i}}{\sigma_{y}^{2}}}_{\Sigma_{w}} = \left(\mathbf{A}\sigma_{y}^{2} + \sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathsf{T}}\right)^{-1} \sum_{i} y_{i} \mathbf{x}_{i}$$

$$\mathbf{w}^{\mathrm{ML}} = \widehat{\mathbf{W}}^{\mathrm{T}} = \left(\sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathrm{T}}\right)^{-1} \sum_{i} y_{i} \mathbf{x}_{i}$$

- The prior acts to "inflate" the apparent covariance of inputs.
- As A is positive (semi)definite, shrinks the weights towards the prior mean (here 0).
- If $A = \alpha I$ this is known as the ridge regression estimator.
- The MAP/shrinkage/ridge weight estimate often has lower squared error (despite bias) and makes more accurate predictions on test inputs than the ML estimate.

As the posterior is Gaussian, the MAP and posterior mean weights are the same:

$$\mathbf{w}^{\text{MAP}} = \underbrace{\left(A + \frac{\sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\text{T}}}{\sigma_{y}^{2}}\right)^{-1}}_{\Sigma_{w}} \underbrace{\frac{\sum_{i} y_{i} \mathbf{x}_{i}}{\sigma_{y}^{2}}}_{\Sigma_{w}} = \left(A\sigma_{y}^{2} + \sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\text{T}}\right)^{-1} \sum_{i} y_{i} \mathbf{x}_{i}$$

$$\mathbf{w}^{\mathrm{ML}} = \widehat{\mathbf{W}}^{\mathrm{T}} = \left(\sum_{i} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathrm{T}}\right)^{-1} \sum_{i} y_{i} \mathbf{x}_{i}$$

- The prior acts to "inflate" the apparent covariance of inputs.
- As A is positive (semi)definite, shrinks the weights towards the prior mean (here 0).
- If $A = \alpha I$ this is known as the ridge regression estimator.
- The MAP/shrinkage/ridge weight estimate often has lower squared error (despite bias) and makes more accurate predictions on test inputs than the ML estimate.
- An example of prior-based regularisation of estimates.

• Models the conditional $P(\mathbf{y}|\mathbf{x})$.

- Models the conditional $P(\mathbf{y}|\mathbf{x})$.
- If we also model $P(\mathbf{x})$, then learning is indistinguishable from unsupervised. In particular if $P(\mathbf{x})$ is Gaussian, and $P(\mathbf{y}|\mathbf{x})$ is linear-Gaussian, then \mathbf{x}, \mathbf{y} are jointly Gaussian.

- Models the conditional $P(\mathbf{y}|\mathbf{x})$.
- ► If we also model $P(\mathbf{x})$, then learning is indistinguishable from unsupervised. In particular if $P(\mathbf{x})$ is Gaussian, and $P(\mathbf{y}|\mathbf{x})$ is linear-Gaussian, then \mathbf{x} , \mathbf{y} are jointly Gaussian.
- Generalised Linear Models (GLMs) generalise to non-Gaussian, exponential-family distributions and to non-linear link functions.

$$egin{aligned} & y_i \sim \mathsf{ExpFam}(\mu_i, \phi) \ & g(\mu_i) = \mathbf{w}^{\mathsf{T}} \mathbf{x}_i \end{aligned}$$

Posterior, or even ML, estimation is not possible in closed form \Rightarrow iterative methods such as gradient ascent or iteratively re-weighted least squares (IRLS). A warning to fMRIers: SPM uses GLM for "general" (not -ised) linear model; which is just linear.

- Models the conditional $P(\mathbf{y}|\mathbf{x})$.
- If we also model $P(\mathbf{x})$, then learning is indistinguishable from unsupervised. In particular if $P(\mathbf{x})$ is Gaussian, and $P(\mathbf{y}|\mathbf{x})$ is linear-Gaussian, then \mathbf{x}, \mathbf{y} are jointly Gaussian.
- Generalised Linear Models (GLMs) generalise to non-Gaussian, exponential-family distributions and to non-linear link functions.

 $y_i \sim \mathsf{ExpFam}(\mu_i, \phi)$ $g(\mu_i) = \mathbf{w}^{\mathsf{T}} \mathbf{x}_i$

Posterior, or even ML, estimation is not possible in closed form \Rightarrow iterative methods such as gradient ascent or iteratively re-weighted least squares (IRLS). A warning to fMRIers: SPM uses GLM for "general" (not -ised) linear model; which is just linear.

These models: Gaussians, Linear-Gaussian Regression and GLMs are important building blocks for the more sophisticated models we will develop later.

- Models the conditional $P(\mathbf{y}|\mathbf{x})$.
- If we also model $P(\mathbf{x})$, then learning is indistinguishable from unsupervised. In particular if $P(\mathbf{x})$ is Gaussian, and $P(\mathbf{y}|\mathbf{x})$ is linear-Gaussian, then \mathbf{x}, \mathbf{y} are jointly Gaussian.
- Generalised Linear Models (GLMs) generalise to non-Gaussian, exponential-family distributions and to non-linear link functions.

 $y_i \sim \mathsf{ExpFam}(\mu_i, \phi)$ $g(\mu_i) = \mathbf{w}^{\mathsf{T}} \mathbf{x}_i$

Posterior, or even ML, estimation is not possible in closed form \Rightarrow iterative methods such as gradient ascent or iteratively re-weighted least squares (IRLS). A warning to fMRIers: SPM uses GLM for "general" (not -ised) linear model; which is just linear.

- These models: Gaussians, Linear-Gaussian Regression and GLMs are important building blocks for the more sophisticated models we will develop later.
- Gaussian models are also used for regression in Gaussian Process Models. We'll see these later too.

What about higher order statistical structure in the data?

What happens if there are outliers?

• There are D(D + 1)/2 parameters in the multivariate Gaussian model. What if D is very large?

What about higher order statistical structure in the data?

 \Rightarrow nonlinear and hierarchical models

What happens if there are outliers?

• There are D(D + 1)/2 parameters in the multivariate Gaussian model. What if D is very large?

What about higher order statistical structure in the data?

 \Rightarrow nonlinear and hierarchical models

What happens if there are outliers?

 \Rightarrow other noise models

• There are D(D + 1)/2 parameters in the multivariate Gaussian model. What if D is very large?

What about higher order statistical structure in the data?

 \Rightarrow nonlinear and hierarchical models

What happens if there are outliers?

 \Rightarrow other noise models

There are D(D + 1)/2 parameters in the multivariate Gaussian model. What if D is very large?

 \Rightarrow dimensionality reduction
End Notes

- It is very important that you understand all the material in the following cribsheet: http://www.gatsby.ucl.ac.uk/teaching/courses/ml1-2014/cribsheet.pdf
- The following notes by (the late) Sam Roweis are quite useful:
 - Matrix identities and matrix derivatives: http://www.cs.nyu.edu/~roweis/notes/matrixid.pdf
 - Gaussian identities: http://www.cs.nyu.edu/~roweis/notes/gaussid.pdf
- Here is a useful statistics / pattern recognition glossary: http://alumni.media.mit.edu/~tpminka/statlearn/glossary/
- Tom Minka's in-depth notes on matrix algebra: http://research.microsoft.com/en-us/um/people/minka/papers/matrix/