

Module MLMI1: Introduction to Machine Learning
 Solutions to Example Sheet 1: Introductory Inference Problems,
 Bayesian Decision Theory, Regression and Classification

Straightforward questions are marked †

*Hard questions are marked **

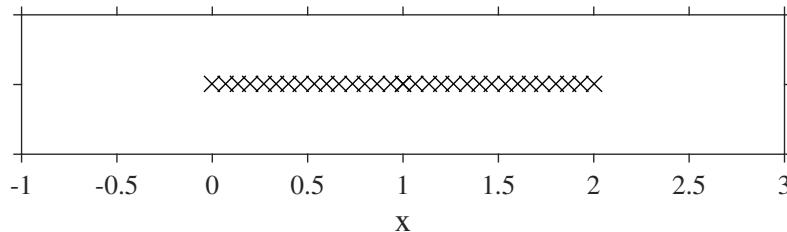
Introductory Inference Problems

1. Maximum likelihood fitting of a Gaussian

- (a) Explain the terms likelihood function, prior probability distribution, and posterior probability distribution, in the context of the inference of parameters θ from data \mathcal{D} .
- (b) A random variable x is believed to have a probability distribution which is Gaussian with mean μ and standard deviation equal to 1,

$$p(x|\mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x - \mu)^2\right).$$

A sample of $N = 32$ data points is collected $\{x_n\}_{n=1}^N$ that are believed to be drawn independently from this distribution. The dataset is shown below:



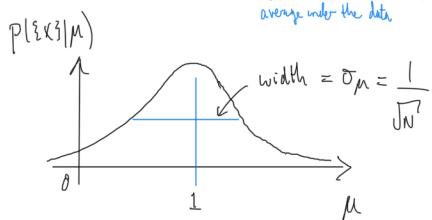
The first and second moments of these data are $\frac{1}{N} \sum_{n=1}^N x_n = 1$ and $\frac{1}{N} \sum_{n=1}^N x_n^2 = 1.3$.

Sketch the likelihood as a function of μ for the dataset. Label the position of the maximum and its width. You do not need to compute the value of the likelihood at its maximum.

$$1. \quad p(x_n | \mu) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x_n - \mu)^2}$$

$$\begin{aligned}
 p(\{x_n\}_{n=1}^N | \mu) &= \prod_{n=1}^N p(x_n | \mu) = \left(\frac{1}{\sqrt{2\pi}} \right)^N e^{-\frac{1}{2} \sum_{n=1}^N (x_n - \mu)^2} \\
 &= (2\pi)^{-N/2} e^{-\frac{1}{2} \left(\sum_{n=1}^N x_n^2 - 2\mu \sum_{n=1}^N x_n + N\mu^2 \right)} \\
 &= (2\pi)^{-N/2} e^{-\frac{1}{2} N \left(\sum_{i=1}^N x_i^2 - 2\mu \sum_{i=1}^N x_i + \mu^2 \right)} \quad \stackrel{\Rightarrow \text{ Gaussian form again}}{\downarrow} \\
 &\approx \mathcal{N}(\mu; \mu_m, \sigma_m^2) = \frac{1}{\sqrt{2\pi \sigma_m^2}} e^{-\frac{1}{2\sigma_m^2} (\mu - \mu_m)^2} \\
 &\approx \frac{1}{\sqrt{2\pi \sigma_m^2}} e^{-\frac{1}{2\sigma_m^2} (\mu^2 - 2\mu \mu_m - \mu_m^2)} \quad \stackrel{\text{compare}}{\downarrow} \quad \stackrel{\text{[Complete the square]}}{\downarrow}
 \end{aligned}$$

$$\Rightarrow \sigma_m^2 = \frac{1}{N} \quad \mu_m = \langle x \rangle = 1 \quad (\text{could also find } \sigma, \text{ but question does not ask for it})$$



N.B. When fitting a Gaussian, the likelihood only depends on the data's 1st & 2nd moments.
So even though the data appear to come from a uniform density here, we only need the two moments.

2. Inference in a Gaussian model

A noisy depth sensor measures the distance to an object an unknown distance d metres away. The depth can be assumed, *a priori*, to be distributed according to a standard Gaussian distribution $p(d) = \mathcal{N}(d; 0, 1)$. The depth sensor returns y a noisy measurement of the depth, that is also assumed to be Gaussian $p(y|d, \sigma_y^2) = \mathcal{N}(y; d, \sigma_y^2)$.

- (a) Compute the posterior distribution over the depth given the observation, $p(d|y, \sigma_y^2)$.
- (b) What happens to the posterior distribution as the measurement noise becomes very large $\sigma_y^2 \rightarrow \infty$? Comment on this result.

The formula for the probability density of a Gaussian distribution of mean μ and variance σ^2 is given by

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

2. $p(d) = \mathcal{N}(d; 0, 1)$ $p(y|d, \sigma_y^2) = \mathcal{N}(y; d, \sigma_y^2)$

$$\begin{aligned} a) \quad p(d|y) &\propto p(y|d, \sigma_y^2) p(d) \propto e^{-\frac{1}{2\sigma_y^2}(y-d)^2 - \frac{1}{2}d^2} \\ &= e^{-\frac{1}{2}\left[\frac{d^2}{\sigma_y^2} + 1\right] - \frac{2dy}{\sigma_y^2} + \frac{y^2}{\sigma_y^2}} \stackrel{\text{Gaussian form}}{\stackrel{\text{conjugate coefficients}}{=}} \\ &= e^{-\frac{1}{2\sigma_{d|y}^2}(d - \mu_{d|y})^2} = e^{-\frac{1}{2}\left(\frac{d^2}{\sigma_{d|y}^2} - \frac{2d\mu_{d|y}}{\sigma_{d|y}^2} + \frac{\mu_{d|y}^2}{\sigma_{d|y}^2}\right)} \end{aligned}$$

$$\therefore \sigma_{d|y}^2 = \frac{1}{\frac{1}{\sigma_y^2} + 1} = \frac{\sigma_y^2}{1 + \sigma_y^2}$$

$$\mu_{d|y} = \sigma_{d|y}^2 \frac{y}{\sigma_y^2} = \frac{y}{1 + \sigma_y^2}$$

$$b) \quad \sigma_y^2 \rightarrow \infty \Rightarrow \sigma_{d|y}^2 \rightarrow 1 \quad (\text{This is the prior variance as it should be as now the sensor is so noisy it does not tell us anything over & above our a priori beliefs})$$

$$\Rightarrow \mu_{d|y} \rightarrow 0 \quad (\text{again collapses back to the prior})$$

Note also that when $\sigma_y^2 \rightarrow 0$ the sensor gives no perfect information about y so $\sigma_{d|y}^2 \rightarrow 0$ & $\mu_{d|y} \rightarrow y$ as expected

3. Bayesian inference for a biased coin*

A sequence of coin tosses are observed from a biased coin $x_{1:N} = \{0, 1, 1, 0, 1, 1, 1, 1, 0\}$ where $x_n = 1$ indicates flip n was a head and $x_n = 0$ indicates that it was tails. An experimenter would like to estimate the coin's probability of landing heads, ρ , from these data.

The experimenter assumes that the coin flips are drawn independently from a Bernoulli distribution $p(x_n|\rho) = \rho^{x_n}(1 - \rho)^{1-x_n}$ and uses a prior distribution of the form

$$p(\rho|n_0, N_0) = \frac{1}{Z(n_0, N_0)} \rho^{n_0} (1 - \rho)^{N_0 - n_0}.$$

Here n_0 and N_0 are parameters set by the experimenter to encapsulate their prior beliefs. $Z(n_0, N_0)$ returns the normalising constant of the distribution as a function of the parameters, n_0 and N_0 .

- (a) Compute the posterior distribution over the bias $p(\rho|x_{1:N}, n_0, N_0)$.
- (b) Compute the *maximum a posteriori* (MAP) estimate for the bias.
- (c) Provide an intuitive interpretation for the parameters of the prior distribution, n_0 and N_0 . For what setting of n_0 and N_0 does the MAP estimate become equal to the maximum-likelihood estimate?

$$\begin{aligned}
 3. a) \quad p(p | x_{1:N}, n_0, N_0) &\propto p(p | n_0, N_0) \prod_{n=1}^N p(x_n | p) \\
 &= \frac{1}{Z(n_0, N_0)} p^{n_0} (1-p)^{N_0 - n_0} p^{\sum_n x_n} (1-p)^{N - \sum_n x_n} \\
 &= \frac{1}{Z} p^{n_0 + \lambda} (1-p)^{N_0 + N - n_0 - \lambda} \quad \text{where } \lambda = \sum_n x_n \\
 &\quad \text{i.e. \# of 1's in dataset}
 \end{aligned}$$

$$\therefore p(p | x_{1:N}, n_0, N_0) = \frac{1}{Z(n^*, N^*)} p^{n^*} (1-p)^{N^* - n^*} \quad \text{where } n^* = n_0 + \lambda \\
 N^* = N_0 + N$$

(This is a Beta distribution & it is conjugate to the likelihood, meaning the posterior has the same form)

$$b) \log p(p | x_{1:N}, n_0, N_0) = -\log Z + n^* \log p + (n^* - n^*) \log (1-p)$$

$$\begin{aligned}
 \frac{d}{dp} \log p(p_{\text{MAP}} | x_{1:N}, n_0, N_0) &= \frac{n^*}{p_{\text{MAP}}} - \frac{(n^* - n^*)}{1-p_{\text{MAP}}} = 0 \\
 (1-p_{\text{MAP}})^{n^*} - p_{\text{MAP}}^{(n^*-n^*)} &= 0 \quad \Rightarrow p_{\text{MAP}} = \frac{n^*}{N^*}
 \end{aligned}$$

c) $N_0 = \# \text{ of pseudo data points seen before real data}$

$n_0 = \# \text{ of 1's in pseudo data}$

ML estimate is recovered when $N_0 = n_0 = 0 \Rightarrow$ flat prior distribution
(no pseudo data)

4. Inferential game show*

On a game show, a contestant is told the rules as follows:

There are four doors, labelled 1, 2, 3 and 4. A single prize has been hidden behind one of them. You get to select one door. Initially your chosen door will not be opened. Instead, the gameshow host will open one of the other three doors, and *he will do so in such a way as not to reveal the prize*. For example, if you first choose door 1, he will then open one of doors 2, 3 and 4, and it is guaranteed that he will choose which one to open so that the prize will not be revealed.

At this point, you will be given a fresh choice of door: you can either stick with your first choice, or you can switch to one of the other closed doors. All the doors will then be opened and you will receive whatever is behind your final choice of door.

- (a) Imagine that the contestant chooses door 1 first; then the gameshow host opens door 4, revealing nothing behind the door, as promised. Should the contestant (a) stick with door 1, or (b) switch to door 2 or 3, or (c) does it make no difference?
- (b) Use Bayes' rule to solve the problem.

4) VLC let door 1 be the one selected by the contestant

let $S = \text{position of prize}$ $S \in \{1, 2, 3, 4\}$

A priori we assume $p(S=k) = 1/4$

The datum we receive after choosing door 1 is either $D=2, D=3, D=4$
i.e. doors 2, 3 or 4 are opened.

Assume that when the host has a choice about which door to open he selects uniformly between those doors not associated with the prize.

i.e.

$$\begin{aligned} p(D=2|S=1) &= 1/3 & p(D=2|S=2) &= 0 & p(D=2|S=3) &= 1/2 & p(D=2|S=4) &= 1/2 \\ p(D=3|S=1) &= 1/3 & p(D=3|S=2) &= 1/2 & p(D=3|S=3) &= 0 & p(D=3|S=4) &= 1/2 \\ p(D=4|S=1) &= 1/3 & p(D=4|S=2) &= 1/2 & p(D=4|S=3) &= 1/2 & p(D=4|S=4) &= 0 \end{aligned}$$

Now apply Bayes' theorem:

$$p(S=k | D=4) = \frac{p(D=4 | S=k) p(S=k)}{p(D=4)}$$

$$\begin{aligned} p(S=1 | D=4) &= \frac{1/3 \cdot 1/4}{p(D=4)} & p(S=2 | D=4) &= \frac{1/2 \cdot 1/4}{p(D=4)} & p(S=3 | D=4) &= \frac{1/2 \cdot 1/4}{p(D=4)} & p(S=4 | D=4) &= 0 \\ &= 1/4 \leftarrow \text{same as prior} & & & & & & \\ & & & = 3/8 \leftarrow \text{greater than prior} & & & = 3/8 \leftarrow \text{greater than prior} & \end{aligned}$$

So, if we switch to doors 2 or 3 we will increase our chances of winning from $1/4$ to $3/8$ i.e. 1.5x.

To get an intuition for the fact that the opening of the door by the host provides information, consider 100 doors & the host opening 98 of them.

This is a version of the Monty Hall problem (see pg 57 of David Mackay's information theory & inference book)

5. Bayesian decision theory*

A data-scientist has computed a complex posterior distribution over a variable of interest, x , given observed data y , that is $p(x|y)$. They would like to return a point estimate of x to their client. The client provides the data-scientist with a reward function $R(\hat{x}, x)$ that indicates their satisfaction with a point estimate \hat{x} when the true state of the variable is x .

- (a) Explain how to use *Bayesian Decision Theory* to determine the optimal point estimate, \hat{x} .
- (b) Compute the optimal point estimate \hat{x} in the case when the reward function is the negative square error between the point estimate and the true value, $R(\hat{x}, x) = -(\hat{x} - x)^2$. Comment on your result.
- (c) Compute the optimal point estimate \hat{x} in the case when the reward function is the negative absolute error between the point estimate and the true value, $R(\hat{x}, x) = -|\hat{x} - x|$. Comment on your result.

5)

$$a) \quad \hat{x}_* = \arg \max_{\hat{x}} \int R(\hat{x}, x) p(x|y) dx$$

$$b) \text{ Find optimum: } - \frac{d}{d\hat{x}} \int (x - \hat{x})^2 p(x|y) dx = 0$$

$$\therefore \int x p(x|y) dx = \hat{x}_*$$

i.e. the posterior mean minimizes the expected squared error

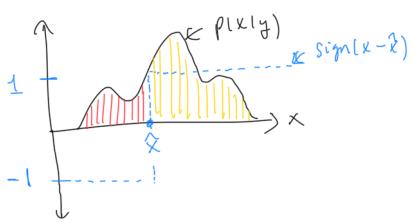
$$c) \text{ Find optimum: } - \frac{d}{d\hat{x}} \int |x - \hat{x}| p(x|y) dx = 0$$

$$= - \frac{d}{d\hat{x}} \int \sqrt{|x - \hat{x}|^2} p(x|y) dx$$

$$= + \frac{1}{2} \cdot 2 \cdot \int \left(\frac{|x - \hat{x}|}{\sqrt{|x - \hat{x}|^2}} \right) p(x|y) dx$$

sign($x - \hat{x}$)

Schematic
of integral



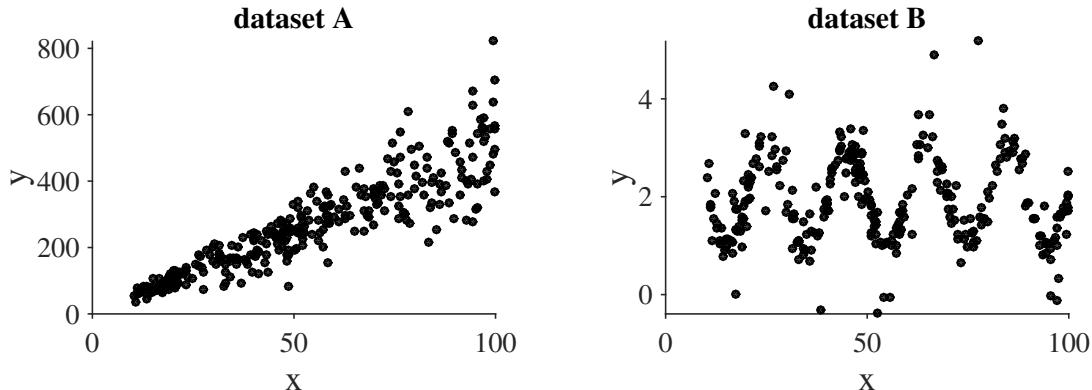
\Rightarrow need to find the point where the
red & yellow areas are equal

\Rightarrow median of the distribution
(has half the density above it
& half below)

Regression

6. Probabilistic models for regression*

A machine learner observes two separate regression datasets comprising scalar inputs and outputs $\{x_n, y_n\}_{n=1}^N$ shown below.



- (a) Suggest a suitable regression model, $p(y_n|x_n)$ for the dataset A. Indicate sensible settings for the parameters in your proposed model where possible. Explain your modelling choices.
- (b) Suggest a suitable regression model, $p(y_n|x_n)$ for the dataset B. Indicate sensible settings for the parameters in your proposed model where possible. Explain your modelling choices.

a) Linear trend, rough gradient ≈ 5 , intercept @ $\{0,0\}$

Noise appears Gaussian but standard deviation grows with x

$$\therefore \text{suggest } \hat{y}_n(x) = 5x + \sigma(x)\varepsilon_n \quad \varepsilon_n \sim N(0,1)$$

$$\text{where } \sigma(x) = |x|$$

Many reasonable choices here, might be good to discuss what people have come up with in the supervision.

b) Sinusoidal trend, time period of rough 25 time steps heavy tailed noise (outliers)

$$\hat{y}_n(x) = 2 + \begin{cases} \sin\left(\frac{2\pi}{25}x\right) & \text{amplitude is } \approx 1 \\ 1 & \text{note mean is 2} \end{cases} + \varepsilon_n \quad \varepsilon_n \sim \text{Student-t}$$

heavy tailed, mean 0
 variance ≈ 1
 (not to guess)
 degree of freedom parameter
 ≈ 2.1 (lower heavier to guess)

Again it's a good one to discuss.

7. Maximum-likelihood learning for a simple regression model

Consider a regression problem where the data comprise N scalar inputs and outputs, $\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\}$, and the goal is to predict y from x .

Assume a very simple linear model, $y_n = ax_n + \epsilon_n$, where the noise ϵ_n is Gaussian with zero mean and variance 1.

- (a) Provide an expression for the log-likelihood of the parameter a .
- (b) Compute the maximum likelihood estimate for a .

$$\begin{aligned}
 a) \quad p(\{y_n\}_{n=1}^N | a, \{x_n\}_{n=1}^N) &= \prod_{n=1}^N p(y_n | a, x_n) \\
 \therefore \log p(\{y_n\}_{n=1}^N | a, \{x_n\}_{n=1}^N) &= \sum_n \left[-\frac{1}{2} \log 2\pi - \frac{1}{2} (y_n - ax_n)^2 \right] \\
 &= -\frac{N}{2} \log 2\pi - \frac{1}{2} \sum_n (y_n - ax_n)^2 = L(a)
 \end{aligned}$$

$$b) \quad \left. \frac{dL(a)}{da} \right|_{a_{ML}} = \sum_n x_n (y_n - a x_n) = 0$$

$$\Rightarrow a_{ML} = \frac{\sum_n x_n y_n}{\sum_n x_n^2}$$

8. Maximum-likelihood learning for multi-output regression*

A data-scientist has collected a regression dataset comprising N scalar inputs ($\{x_n\}_{n=1}^N$) and N scalar outputs ($\{y_n\}_{n=1}^N$). Their goal is to predict y from x and they have assumed a very simple linear model, $y_n = ax_n + \epsilon_n$.

The data-scientist also has access to a second set of outputs ($\{z_n\}_{n=1}^N$) that are well described by the model $z_n = x_n + \epsilon'_n$.

The noise variables ϵ_n and ϵ'_n are known to be zero mean correlated Gaussian variables

$$p\left(\begin{bmatrix} \epsilon_n \\ \epsilon'_n \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \epsilon_n \\ \epsilon'_n \end{bmatrix}; \mathbf{0}, \Sigma\right) \text{ where } \Sigma^{-1} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}.$$

- (a) Provide an expression for the log-likelihood of the parameter a .
- (b) Compute the maximum likelihood estimate for a .
- (c) Do the additional outputs $\{z_n\}_{n=1}^N$ provide useful additional information for estimating a ? Explain your reasoning.

The formula for the probability density of a multivariate Gaussian distribution of mean μ and covariance Σ is given by

$$\mathcal{N}(\mathbf{x}; \mu, \Sigma) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^\top \Sigma^{-1} (\mathbf{x} - \mu)\right).$$

$$\begin{aligned}
 a) L(a) &= \log p(y_1, z_1, \dots, y_N, z_N | x_1, \dots, x_N, a) = \sum_n \log p(y_n, z_n | x_n, a) \\
 \text{where } p(y_n, z_n | x_n, a) &= N\left(\begin{bmatrix} y_n \\ z_n \end{bmatrix} ; \begin{bmatrix} ax_n \\ x_n \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}^{-1}\right) \\
 \therefore L(a) &= \sum_n -\frac{1}{2} \left(\begin{bmatrix} y_n \\ z_n \end{bmatrix} - \begin{bmatrix} ax_n \\ x_n \end{bmatrix} \right)^\top \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix} \left(\begin{bmatrix} y_n \\ z_n \end{bmatrix} - \begin{bmatrix} ax_n \\ x_n \end{bmatrix} \right) - \frac{N}{2} \log(2\pi|\Sigma|) \\
 &= -\frac{1}{2} \sum_n (y_n - ax_n)^2 - \frac{1}{2} \sum_n (y_n - ax_n)(z_n - x_n) - \frac{1}{2} \sum_n (z_n - x_n)^2 - \frac{N}{2} \log(2\pi|\Sigma|) \\
 b) \frac{dL(a)}{da} &= \sum_n (y_n - ax_n)x_n + \frac{1}{2} \sum_n (z_n - x_n)x_n = 0 \\
 &= \sum_n y_n x_n + \frac{1}{2} \sum_n z_n x_n - \frac{1}{2} \sum_n x_n^2 - a \sum_n x_n^2 \\
 \therefore a &= \left(\sum_n y_n x_n + \frac{1}{2} \sum_n (z_n - x_n)x_n \right) / \sum_n x_n^2 \quad (\text{max likelihood estimate}) \\
 &\quad \text{bit from just observing } y_n \quad \text{extra bit from observing } z_n \\
 &\quad \text{new contribution from observing } z_n
 \end{aligned}$$

- c) The additional outputs change the ML estimate of a . This means that they must provide useful information about a . They do this because the noise in z_n is correlated with the noise in y_n & so observing z_n reveals information about the noise ϵ_n & allows more accurate identification of a .

9. Bayesian linear regression.

A single data point $\{x, y\}$ is fit using Bayesian linear regression. The output y is assumed to be generated from the input x according to a linear relationship that is corrupted by Gaussian noise $y = mx + c + \epsilon$. The noise ϵ is mean 0 and variance 1 so $p(y|m, c, x) = \mathcal{N}(y; mx + c, 1)$. Gaussian priors are placed on the slope m and intercept c with zero mean and unit variance, that is $p(m) = \mathcal{N}(m; 0, 1)$ and $p(c) = \mathcal{N}(c; 0, 1)$.

- (a) Compute the posterior probability of the slope and intercept given the data point, that is $p(m, c|x, y)$.

$$a) p(m, c | y, x) = \frac{p(y|m, c, x) p(m)p(c)}{p(y|x)}$$

$$\propto p(y|m, c, x) p(m)p(c)$$

$= N\left(\begin{bmatrix} m \\ c \end{bmatrix} ; \underline{\mu}_{\text{post}}, \underline{\Sigma}_{\text{post}}\right)$

tactic ← product of Gaussian density in M&C

① Substitute in for these densities and ② identify mean & covariance by comparing coefficients

② ← posterior must be Gaussian: just need to compute mean & covariance

$$\begin{aligned} ① p(m, c | y, x, \sigma_y^2) &\propto \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \frac{(y - mx - c)^2}{\sigma_y^2}} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \frac{m^2}{\sigma_y^2}} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \frac{c^2}{\sigma_y^2}} \\ &\propto e^{-\frac{1}{2} \frac{y^2}{\sigma_y^2} - \frac{1}{2} \frac{m^2}{\sigma_y^2} \left(x^2 + 1 \right) - \frac{1}{2} \frac{c^2}{\sigma_y^2} - \frac{1}{2} \frac{mx + y}{\sigma_y^2}} \\ &\propto e^{-\frac{1}{2} \left(\begin{bmatrix} m \\ c \end{bmatrix} - \underline{\mu}_{\text{post}} \right) \left(\underline{\Sigma}_{\text{post}}^{-1} \left(\begin{bmatrix} m \\ c \end{bmatrix} - \underline{\mu}_{\text{post}} \right) \right)} \\ &\propto e^{-\frac{1}{2} \left[\begin{bmatrix} m \\ c \end{bmatrix}^T \underline{\Sigma}_{\text{post}}^{-1} \begin{bmatrix} m \\ c \end{bmatrix} + \begin{bmatrix} m \\ c \end{bmatrix}^T \underline{\Sigma}_{\text{post}}^{-1} \underline{\mu}_{\text{post}} \right]} \\ &\Rightarrow \underline{\Sigma}_{\text{post}} = \begin{bmatrix} x^2 + 1 & x \\ x & 2 \end{bmatrix}^{-1} = \frac{1}{(x^2 + 1)(2) - x^2} \begin{pmatrix} 2 & -x \\ -x & x^2 + 1 \end{pmatrix} \\ &= \frac{1}{x^2 + 2} \begin{pmatrix} 2 & -x \\ -x & x^2 + 1 \end{pmatrix} \end{aligned}$$

inverse covariance times mean

$$\& \text{since } \left(\underline{\Sigma}_{\text{post}} \right)^{-1} \underline{\mu}_{\text{post}} = \begin{bmatrix} yx \\ y \end{bmatrix} \quad \underline{\Sigma}_{\text{post}}$$

$$\underline{\mu}_{\text{post}} = \frac{1}{x^2 + 2} \begin{pmatrix} 2 & -x \\ -x & x^2 + 1 \end{pmatrix} \begin{bmatrix} yx \\ y \end{bmatrix}$$

$$= \frac{1}{x^2 + 2} \begin{bmatrix} 2yx - yx \\ -yx^2 + (x^2 + 1)y \end{bmatrix} = \frac{1}{x^2 + 2} \begin{bmatrix} yx \\ y \end{bmatrix}$$

- (b) Show that the posterior derived in part a is consistent with the expressions for Bayesian linear regression given in lectures.

b) Check these expressions match with those in lectures : $\underline{y} = \underline{\tilde{x}} \underline{w} + \underline{\varepsilon}$ $\underline{\varepsilon} \sim N(0, \sigma^2 \underline{\underline{I}})$

$$p(\underline{w} | \text{Data}, \sigma^2, \lambda) = N(\underline{w}; \underline{\mu}_{\text{WID}}, \underline{\Sigma}_{\text{WID}}) \quad \underline{w} \sim N(0, \lambda^{-1} \underline{\underline{I}})$$

$$\text{where } \underline{\Sigma}_{\text{WID}} = \left(\frac{1}{\lambda} \underline{\underline{I}} + \frac{1}{\sigma^2} \underline{\tilde{x}}^T \underline{\tilde{x}} \right)^{-1} \quad \underline{\mu}_{\text{WID}} = \underline{\Sigma}_{\text{WID}} \frac{1}{\sigma^2} \underline{\tilde{x}}^T \underline{y}$$

in our case $\sigma^2 = \lambda = 1$ (prior & observation noise are unit variance)

$$\underline{\tilde{x}} = [x, 1]$$

$$\underline{w} = \begin{bmatrix} m \\ c \end{bmatrix}$$

$$\Rightarrow \underline{\Sigma}_{\text{WID}} = \left(\underline{\underline{I}} + \begin{bmatrix} x \\ 1 \end{bmatrix} \begin{bmatrix} x & 1 \end{bmatrix} \right)^{-1} = \begin{pmatrix} x^2 + 1 & x \\ x & 2 \end{pmatrix}^{-1} \quad \checkmark$$

$$\underline{\mu}_{\text{WID}} = \begin{pmatrix} x^2 + 1 & x \\ x & 2 \end{pmatrix}^{-1} \begin{bmatrix} x \\ 1 \end{bmatrix} \underline{y} \quad \checkmark$$

(c) Compute the posterior in the following three cases and provide explanations for why the posteriors take the form that they do.

- i. $x = 0$ and $y = 0$
- ii. $x = 1$ and $y = 0$
- iii. $x = 100$ and $y = 100$

c) plug the following settings $\begin{matrix} x = 0 \\ y = 0 \end{matrix}$ $\begin{matrix} x = 1 \\ y = 1 \end{matrix}$ $\begin{matrix} x = 100 \\ y = 100 \end{matrix}$

into :

$$\underline{\Sigma}^{\text{post}} = \frac{1}{x^2+2} \begin{bmatrix} 2 & -x \\ -x & x^2+1 \end{bmatrix}$$

$$\underline{\mu}^{\text{post}} = \frac{1}{x^2+2} \begin{bmatrix} y & x \\ 0 & 0 \end{bmatrix}$$

i) $x = 0 \quad y = 0 \quad \underline{\Sigma}^{\text{post}} = \begin{bmatrix} 1 & 0 \\ 0 & 1/2 \end{bmatrix} \quad \underline{\mu}^{\text{post}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

ii) $x = 1 \quad y = 0 \quad \underline{\Sigma}^{\text{post}} = \begin{bmatrix} 2/3 & -1/3 \\ -1/3 & 2/3 \end{bmatrix} \quad \underline{\mu}^{\text{post}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

iii) $x = 100 \quad y = 100 \quad \underline{\Sigma}^{\text{post}} \approx \begin{bmatrix} 2/100^2 & -1/100 \\ -1/100 & 1 \end{bmatrix} \quad \underline{\mu}^{\text{post}} \approx \begin{bmatrix} 1 \\ 1/100 \end{bmatrix}$

i) data doesn't tell us anything about m : $y = mx + c$ & $x = 0$

This is why posterior mean & variance of m stays @ prior (mean 0 Variance 1)

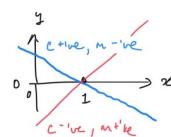
data tells us that c is likely to be ≈ 0 , but the observation noise is quite large

so only weak evidence for this. This why posterior mean is 0 and posterior variance is $1/2$ (prior variance was 1).

ii) The data can be explained by either positive c & negative m or negative c & positive m :

Hence posterior covariance between m & c is negative ($-1/3$)

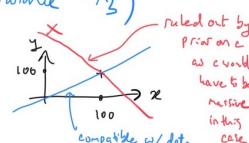
Dataset provides weak evidence that m & c are both small in magnitude & due to symmetry both still have posterior mean 0 (uncertainty has reduced from variance 1 to variance $2/3$)



iii) The data $y = 100 \quad x = 100$ rules out negative m :

We are therefore quite certain that m is close to 1 (posterior mean very close to this value & variance $\approx 2/100^2$)

The data give almost no information about c though with posterior mean & variance close to the prior -



You might like to compute the predictive mean for these three cases i.e. the mean of $p(y^*|x^*, x, y)$ where x^* is the location of a new test input and y^* is the corresponding output.

Predictive mean : use $y^* = Mx^* + C + \varepsilon^*$ now average wrt $p(M, C, \varepsilon^* | x, y)$:

$$\mathbb{E}_{\substack{p(M, C, \varepsilon^* | x, y) \\ \text{independent from } x^* \text{ & } y}} [Mx^* + C + \varepsilon^*] = \mu_M^{\text{post}} x^* + \mu_C^{\text{post}} + b$$

i) predictive mean = 0 (horizontal line)

ii) predictive mean = 0 (horizontal line)

iii) predictive mean = $\frac{1}{100^2 + 2} \left(100^2 x^* + 100 \right)$ (line with gradient v. close to 1 and a small positive intercept)

Classification

10. Probit Classification

Consider classification as described in lectures, but with the following model

$$y^{(n)} = H(\mathbf{w}^\top \mathbf{x}^{(n)} + \epsilon_n)$$

where ϵ_n is Gaussian with mean 0 and variance σ^2 and $H(\cdot)$ is the Heaviside step function.

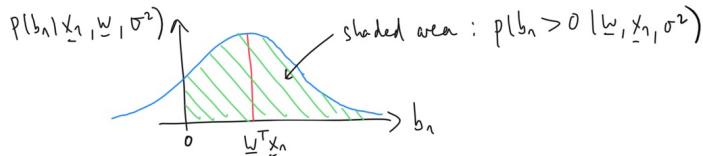
- (a) Compute the probability $P(y^{(n)} = 1 | \mathbf{x}_n, \mathbf{w}, \sigma^2)$ in terms of the Gaussian cumulative distribution. Sketch $P(y^{(n)} = 1 | \mathbf{x}_n, \mathbf{w}, \sigma^2)$ as a function of the inputs \mathbf{x}_n in the case where they are one dimensional.

- (b) What happens as the noise variance tends to infinity $\sigma^2 \rightarrow \infty$?

a) $y_n = H(\underline{\mathbf{w}}^\top \underline{\mathbf{x}}_n + \varepsilon_n) \quad \varepsilon_n \sim N(0, \sigma^2)$

$$\text{let } b_n = \underline{\mathbf{w}}^\top \underline{\mathbf{x}}_n + \varepsilon_n \quad b_n \sim N(\underline{\mathbf{w}}^\top \underline{\mathbf{x}}_n, \sigma^2)$$

$$\therefore P(y_n=1 | \underline{\mathbf{w}}, \underline{\mathbf{x}}_n, \sigma^2) = P(b_n > 0 | \underline{\mathbf{w}}, \underline{\mathbf{x}}_n, \sigma^2)$$



$$\therefore P(y_n=1 | \underline{\mathbf{w}}, \underline{\mathbf{x}}_n, \sigma^2) = \int_0^\infty N(b_n; \underline{\mathbf{w}}^\top \underline{\mathbf{x}}_n, \sigma^2) db_n$$

Let $v_n = -\left(\frac{b_n - \underline{\mathbf{w}}^\top \underline{\mathbf{x}}_n}{\sigma^2}\right)$ & transform to integral over this quantity (v_n will be distributed according to a standard normal)

$$P(y_n=1 | \underline{\mathbf{w}}, \underline{\mathbf{x}}_n, \sigma^2) = \int_{-\infty}^{\frac{\underline{\mathbf{w}}^\top \underline{\mathbf{x}}_n}{\sigma^2}} N(v_n; 0, 1) dv_n$$

$$P(y_n=1 | \underline{\mathbf{w}}, \underline{\mathbf{x}}_n, \sigma^2) = \text{CDF}\left(\frac{\underline{\mathbf{w}}^\top \underline{\mathbf{x}}_n}{\sigma^2}\right)$$

b) $\sigma^2 \rightarrow \infty \Rightarrow P(y_n=1 | \underline{\mathbf{x}}_n, \underline{\mathbf{w}}, \sigma^2) = 1/2$

The noise swamps $\underline{\mathbf{w}}^\top \underline{\mathbf{x}}_n$ resulting in an output which is a coin toss

11. Multi-class Classification

Consider a multi-class classification problem with K classes. The training labels are represented by K dimensional vectors \mathbf{t}_n which has a single element set to 1, indicating the class membership, and all other values are set to 0. The inputs are multi-dimensional vectors \mathbf{x}_n . The goal is to use a training set of input vectors and output labels $\{\mathbf{x}_n, \mathbf{t}_n\}_{n=1}^N$ to enable prediction at unseen input locations.

A friend suggests using a soft-max function for this purpose which is parameterised by weights $\mathbf{W} = \{\mathbf{w}_k\}_{k=1}^K$. The output of the function is a vector, \mathbf{y} , with elements given by

$$y_i(\mathbf{x}; \mathbf{W}) = \frac{\exp(\mathbf{w}_i^\top \mathbf{x})}{\sum_{k=1}^K \exp(\mathbf{w}_k^\top \mathbf{x})}.$$

- (a) What happens to the softmax function as the magnitude of the weights tends to infinity?
- (b) Interpreting the output of the softmax as $y_i = p(t_i = 1 | \mathbf{W}, \mathbf{x})$ write down a cost-function for training this network based on the log-probability of the training data given the weights \mathbf{W} and inputs $\{\mathbf{x}_n\}_{n=1}^N$.
- (c) What is the relationship between this network and logistic regression?

$$1. \quad y_i(x; \underline{w}) = \frac{e^{\underline{w}_i^\top \underline{x}}}{\sum_{k=1}^K e^{\underline{w}_k^\top \underline{x}}} \quad \text{let } \hat{\underline{w}}_{\max} = \underset{\hat{\underline{w}}_i \in \hat{\underline{w}}_1, \dots, \hat{\underline{w}}_K}{\arg \max} \hat{\underline{w}}_i^\top \underline{x}$$

a) let $\underline{w}_i = \beta \hat{\underline{w}}_i$ unit vector in direction of \underline{w}_i
magnitude

$$y_i(x; \underline{w}) = \frac{e^{\beta \hat{\underline{w}}_i^\top \underline{x}}}{\sum_{k=1}^K e^{\beta \hat{\underline{w}}_k^\top \underline{x}}} = \frac{e^{\beta (\hat{\underline{w}}_i - \hat{\underline{w}}_{\max})^\top \underline{x}}}{\sum_{k=1}^K e^{\beta (\hat{\underline{w}}_k - \hat{\underline{w}}_{\max})^\top \underline{x}}}$$

$(\hat{\underline{w}}_i - \hat{\underline{w}}_{\max})^\top \underline{x}$ are a set of K scalars, one of which is zero & the others are negative

$$\therefore \text{as } \beta \rightarrow \infty \quad e^{\beta (\hat{\underline{w}}_i - \hat{\underline{w}}_{\max})^\top \underline{x}} \rightarrow \begin{cases} 0 & \text{if } \hat{\underline{w}}_i \neq \hat{\underline{w}}_{\max} \\ \infty & \text{if } \hat{\underline{w}}_i = \hat{\underline{w}}_{\max} \end{cases}$$

$$\Rightarrow y_i(x; \underline{w}) \rightarrow \begin{cases} 0 & \text{if } i \neq i_{\max} \\ 1 & \text{if } i = i_{\max} \end{cases} \quad \text{hence the general name "softmax"}$$

This is a zero/one encoding of the arg max function:

$$\text{ie as } \beta \rightarrow \infty \quad y_i(x; \underline{w}) \rightarrow \prod_{k=1}^K \left(\underset{k}{\arg \max} \hat{\underline{w}}_k^\top \underline{x} = i \right)$$

NB
many people will be able to intuit this result without going through all of this.

indicator function that takes the value 1 if the argument is true & is otherwise 0.

$$b) \quad y_i = p(t_i=1 | w, x)$$

$$\begin{aligned} P\left(\left\{\underline{t}^{(n)}\right\}_{n=1}^N \mid \underline{w}, \underline{x}^{(n)}\right) &= \prod_{n=1}^N P(t^{(n)} \mid \underline{w}, \underline{x}^{(n)}) \\ &= \prod_{n=1}^N \prod_{k=1}^K y_k(x^{(n)}, w)^{t_k^{(n)}} \\ &= e^{\sum_n \sum_k t_k^{(n)} \log y_k(x^{(n)}, w)} \end{aligned}$$

\Rightarrow cost function could be

$$\underset{\underline{w}}{\text{arg max}} \sum_{n=1}^N \sum_{k=1}^K t_k^{(n)} \log y_k(x^{(n)}, w)$$

c) Consider $K=2$ (two classes)

$$\begin{aligned} y_1(x, w) &= \frac{e^{w^\top x}}{e^{w_1^\top x} + e^{w_2^\top x}} = \frac{1}{1 + e^{(w_2 - w_1)^\top x}} \\ &= \frac{1}{1 + e^{v^\top x}} \Rightarrow \text{equivalent to logistic} \\ &\quad \text{classification when there are 2 classes} \end{aligned}$$