

# PREPARATORY WORK

## DEFINITIONS

- PRIOR PROBABILITY DISTRIBUTION (often called PRIOR) of a certain quantity is the probability distribution that would express one's beliefs about this quantity before some evidence is taken into account.

- LIKELIHOOD FUNCTION measures the goodness fit of a statistical model to a sample of data for a given values of unknown parameters.  $P(X|\theta)$  (prob. of evidence given parameters)

$$\mathcal{L}(\theta) = \mathcal{L}(\theta|x) = P_\theta(x) = P_\theta(X=x) \quad \text{DISCRETE}$$

$$\mathcal{L}(\theta) = \mathcal{L}(\theta|x) = f_{\theta|x}(x) \quad \text{CONTINUOUS}$$

- POSTERIOR PROBABILITY DISTRIBUTION is the probability distribution of an unknown quantity, treated as random variable, conditional on the evidence obtained from an experiment or survey.

$$P(\theta|x) \quad (\text{prob. of parameters given the evidence})$$

Given a prior belief that a probability distribution is  $P(\theta)$  and that the observation  $x$  have a likelihood  $P(x|\theta)$ , then the posterior probability is defined as:

$$P(\theta|x) = \frac{P(x|\theta) \cdot P(\theta)}{P(x)} \quad \text{where } P(x) \text{ is the normalizing const.}$$

$$= \int P(x|\theta) \cdot P(\theta) d\theta \quad \sum_{\theta \in \Theta} P(x|\theta) P(\theta)$$

Obs: the posterior probability is proportional to likelihood • prior prob.

- MAXIMUM LIKELIHOOD ESTIMATION (MLE): method of estimating the parameters of an assumed probability distribution given some observed data. This is achieved by maximizing the likelihood function.

$$\hat{\theta}_{\text{ML}} = \underset{\theta}{\operatorname{argmax}} \mathcal{L}(\theta) = \underset{\theta}{\operatorname{argmax}} \mathcal{L}(\theta | y)$$

↑ notation  
observed data  
maybe:  $\mathcal{L}(\theta | x)$

$$= \underset{\theta}{\operatorname{argmax}} P(x|\theta)$$

- MAXIMUM A POSTERIORI ESTIMATION (MAP): method of estimating the parameters of a statistical model given some data. It's done by maximizing posterior prob. distribution. It's closely related to ML estimation, but employs an augmented optimization objective, which incorporates the prior distribution.

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} P(\theta|x) = \underset{\theta}{\operatorname{argmax}} \frac{P(x|\theta) \cdot P(\theta)}{P(x)}$$

const.  $\rightarrow P(x)$

[4.]

a) done

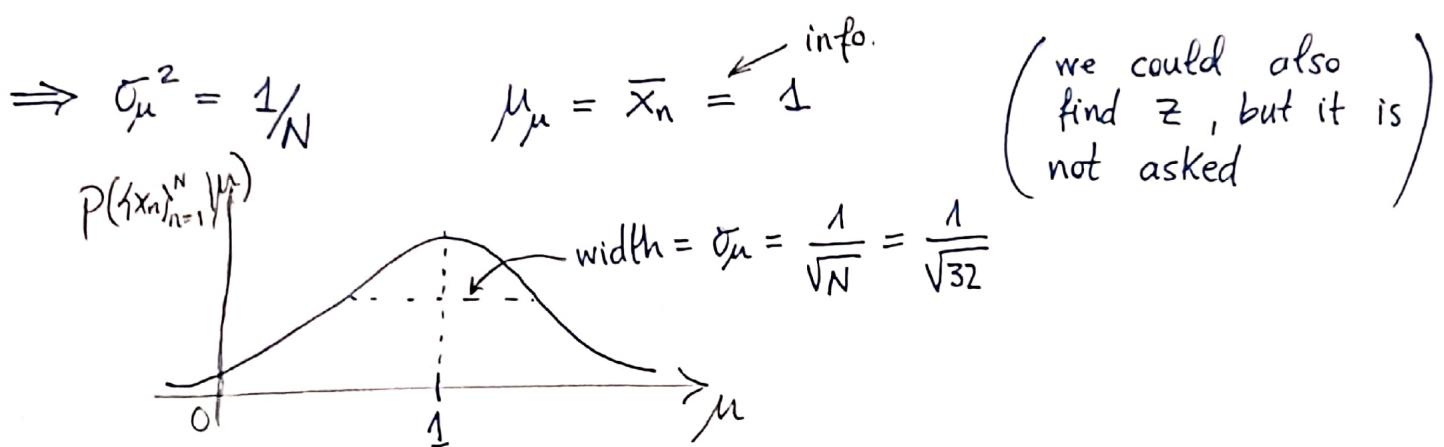
$$b) p(x|\mu) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu)^2}$$

$$\begin{aligned} p(\{x_n\}_{n=1}^N | \mu) &= \prod_{n=1}^N p(x_n | \mu) = \left( \frac{1}{\sqrt{2\pi}} \right)^N e^{-\frac{1}{2} \sum_{n=1}^N (x_n - \mu)^2} = \\ &= (2\pi)^{-\frac{N}{2}} \cdot e^{-\frac{1}{2} \sum_{n=1}^N (x_n^2 - 2x_n\mu + \mu^2)} = \\ &= (2\pi)^{-\frac{N}{2}} \cdot e^{-\frac{1}{2}N(\bar{x}_n^2 - 2\mu\bar{x}_n + \mu^2)} \end{aligned}$$

compare  
coefficients

On the other hand:

$$\begin{aligned} z N(\mu; \mu_\mu, \sigma_\mu^2) &= \frac{z}{\sqrt{2\pi \sigma_\mu^2}} e^{-\frac{1}{2\sigma_\mu^2}(\mu - \mu_\mu)^2} = \\ &= \frac{z}{\sqrt{2\pi \sigma_\mu^2}} e^{\frac{-1}{2\sigma_\mu^2}(\mu^2 - 2\mu\mu_\mu + \mu_\mu^2)} \end{aligned}$$



NB: when fitting a Gaussian, the likelihood only depend on the data's 1st and 2nd moments, so even though the data appear to come from a uniform density here, we only need the two moments.

2. Info:  $p(d) = N(d; 0, 1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}d^2}$   
 $p(y|d, \sigma_y^2) = N(y; d, \sigma_y^2) = \frac{1}{\sqrt{2\pi\sigma_y^2}} e^{-\frac{1}{2\sigma_y^2}(y-d)^2}$

a) Posterior:  $p(d|y) \propto \underbrace{p(y|d, \sigma_y^2)}_{\text{likelihood}} \cdot \underbrace{p(d)}_{\text{prior}}$   
 $\propto e^{-\frac{1}{2\sigma_y^2}(y-d)^2 - \frac{1}{2}d^2} = e^{-\frac{1}{2\sigma_y^2}(y^2 - 2yd + d^2) - \frac{1}{2}d^2} =$   
 $= e^{-\frac{1}{2}\left[d^2\left(\frac{1}{\sigma_y^2} + 1\right) - \frac{2dy}{\sigma_y^2} + \frac{y^2}{\sigma_y^2}\right]} \xleftarrow{\text{compare coefficients}}$

If we suppose that:  $p(d|y) \sim N(d; \mu_{d|y}, \sigma_{d|y}^2)$

$$e^{-\frac{1}{2}\left[\frac{d^2}{\sigma_{d|y}^2} - 2\frac{d\mu_{d|y}}{\sigma_{d|y}^2} + \frac{\mu_{d|y}^2}{\sigma_{d|y}^2}\right]}$$

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

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

b) If  $\sigma_y^2 \rightarrow \infty \Rightarrow \sigma_{d|y}^2 \rightarrow 1$

$\parallel$   
 $\xrightarrow{\quad} \mu_{d|y} \rightarrow 0$

collapses to the prior, so the posterior doesn't give information

this is the prior variance as it should be as now the sensor is so noisy it doesn't tell us anything  
 $\rightarrow$  we just have a prior belief

NB: note also that when  $\sigma_y^2 \rightarrow 0$  the sensor gives us perfect information about  $y$ , so  $\sigma_{d|y}^2 \rightarrow 0$  and  $\mu_{d|y} \rightarrow y$  as expected.

3.  $x_{1:N} = \underbrace{\{0, 1, 1, 1, 0, 1, 0, 0, \dots\}}_N$  coinflips

$x_n = 1$  heads  
 $x_n = 0$  tails  
 $p :=$  prob. of landing heads

$$P(x_n | p) = p^{x_n} (1-p)^{1-x_n}$$

$$P(p | n_0, N_0) = \frac{1}{Z(n_0, N_0)} p^{n_0} (1-p)^{N_0 - n_0}; \quad Z(n_0, N_0) := \text{normalizing constant}$$

a)  $\overbrace{P(p | x_{1:N}, n_0, N_0)}^{\text{posterior}} \propto \overbrace{P(x_{1:N} | p, n_0, N_0) \cdot P(p | n_0, N_0)}^{\text{likelihood prior}} =$

$$= P(p | n_0, N_0) \cdot \prod_{i=1}^N P(x_i | p, n_0, N_0) =$$

$$= \frac{1}{Z(n_0, N_0)} p^{n_0} (1-p)^{N_0 - n_0} \cdot \prod_{i=1}^N p^{x_i} (1-p)^{1-x_i} =$$

$$= \frac{1}{Z(n_0, N_0)} p^{n_0} (1-p)^{N_0 - n_0} \cdot p^{\sum x_i} \cdot (1-p)^{N - \sum x_i} =$$

$$= \frac{1}{Z(n_0, N_0)} p^{n_0+n} (1-p)^{N_0+N-n_0-n} \equiv \left[ \frac{1}{Z(n', N')} p^{n'} (1-p)^{N'-n'} \right]$$

number of 1's  
 $n = \sum x_i$

(this is a Beta distribution and it is conjugate to the likelihood, meaning that the posterior has the same form)

where  $\begin{cases} n' = n_0 + n \\ N' = N_0 + N \end{cases}$

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

$$\Rightarrow \frac{d}{dp} \log P(p | x_{1:N}, n_0, N_0) = \frac{n'}{p_{MAP}} - \frac{N' - n'}{1 - p_{MAP}} = 0 \Rightarrow \boxed{p_{MAP} = \frac{n'}{N'}}$$

c)  $N_0$  = number of pseudo data points, i.e., points seen before real data  
 $n_0$  = number of 1's in pseudo data

ML estimate = MAP estimate when prior is constant, so  
 ML estimate is recovered when  $N_0 = n_0 = 0$  (flat prior distrib.)  
 i.e. no pseudo data

4. WLG 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  $P(S=k) = \frac{1}{4}$   $k=1, \dots, 4$ .

The datum we receive after choosing door 1 is either  $D=2, D=3, D=4$ , i.e., the gameshow host opens door 2, 3 or 4.

We assume that when the host has a choice about which door to open he selects between those doors without prize.

$P(D=2 S=1) = \frac{1}{3}$	$P(D=2 S=2) = 0$	$P(D=2 S=3) = \frac{1}{2}$	$P(D=2 S=4) = \frac{1}{2}$
$P(D=3 S=1) = \frac{1}{3}$	$P(D=3 S=2) = \frac{1}{2}$	$P(D=3 S=3) = 0$	$P(D=3 S=4) = \frac{1}{2}$
$P(D=4 S=1) = \frac{1}{3}$	$P(D=4 S=2) = \frac{1}{2}$	$P(D=4 S=3) = \frac{1}{2}$	$P(D=4 S=4) = 0$

Applying Bayes' theorem:

$$P(S=k|D=4) = \frac{\overbrace{P(D=4|S=k) \cdot P(S=k)}^{=\frac{1}{4} \text{ (prior)}}}{\underbrace{P(D=4)}_{\approx \frac{1}{3}}}$$

$$P(S=1|D=4) = \frac{P(D=4|S=1) P(S=1)}{P(D=4)} = \frac{\frac{1}{3} \cdot \frac{1}{4}}{\frac{1}{3}} = \frac{1}{4} \quad (\text{same as prior})$$

$$P(S=2|D=4) = \frac{\frac{1}{2} \cdot \frac{1}{4}}{\frac{1}{3}} = \frac{3}{8} \quad ; \quad P(S=3|D=4) = \frac{\frac{1}{2} \cdot \frac{1}{4}}{\frac{1}{3}} = \frac{3}{8}$$

$$P(S=4|D=4) = 0$$

greater than prior

$\Rightarrow$  So if we switch to doors 2 or 3 we'll increase our chances of winning from  $\frac{1}{4}$  to  $\frac{3}{8}$  (i.e. 1.5x)

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

NB2: Version of 'Monty Hall problem' (page 97, David McKay, Information Theory, Inference book)

5.

a)  $\hat{x}_* = \underset{\hat{x}}{\operatorname{argmax}} \int R(\hat{x}, x) P(x|y) dx$

b) Find optimum:  $R(\hat{x}, x) = -(\hat{x} - x)^2$

$$\frac{d}{d\hat{x}} \int -(\hat{x} - x)^2 P(x|y) dx = 0 \iff 2 \int (x - \hat{x}) P(x|y) dx = 0$$

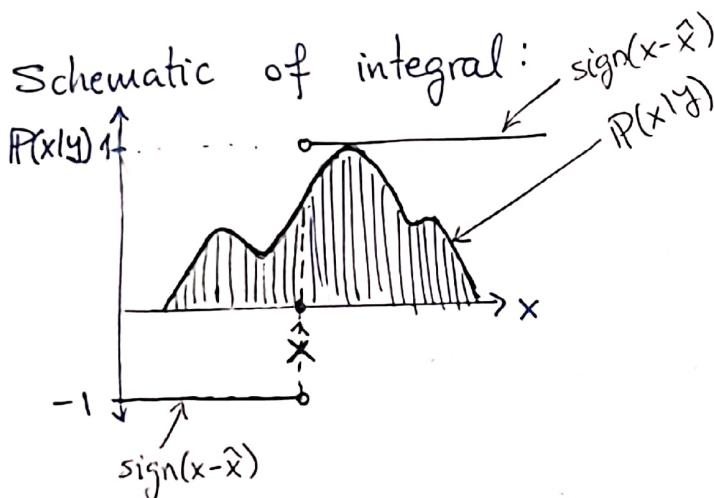
$$\iff \int \hat{x} P(x|y) dx = \underbrace{\int \hat{x} P(x|y) dx}_{=1} = \int x P(x|y) dx \implies \boxed{\hat{x}_* = \int x P(x|y) dx}$$

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

c) Find optimum when  $R(\hat{x}, x) = -|\hat{x} - x|$

$$\frac{d}{d\hat{x}} \int -|\hat{x} - x| P(x|y) dx = 0 \iff -\frac{d}{d\hat{x}} \int \sqrt{(\hat{x} - x)^2} P(x|y) dx = 0$$

$$\iff \frac{1}{2} \cdot 2 \int \underbrace{\frac{1}{\sqrt{(\hat{x} - x)^2}}}_{\text{sign}(x - \hat{x})} P(x|y) dx = 0 \iff \int \text{sign}(x - \hat{x}) P(x|y) dx = 0$$



$\Rightarrow$  need to find the point where the green and blue areas are equal

It is called the median of the distribution,  $\hat{x}_*$ , i.e. it has half the density above it and half below.

6.

- a) LINEAR REGRESSION MODEL : gradient  $\approx 5$  and intercept  $(0,0)$   
 The noise appears to be Gaussian but its standard dev. appears to grow with  $x$ . Therefore an appropriate model (including noise) would be:  $y_n(x) = 5x + \sigma(x)\varepsilon_n$   
 where  $\varepsilon_n \sim N(0,1)$  and  $\sigma(x) = |x|$

Note that there are many reasonable choices that could be made to model this data.

- b) SINUSOIDAL MODEL : period  $\approx 25$  time steps  
 mean = 2  
 amplitude  $\approx 1$

Since there are outliers present, an appropriate noise model would be some heavy tailed model, such as a Student-t distribution:  $y_n(x) = \underbrace{2}_{\text{mean}} + \underbrace{\sin\left(\frac{2\pi}{25}x\right)}_{\substack{\text{amplitude} \\ \text{is } \approx 1}} + \underbrace{\varepsilon_n}_{\substack{\text{time steps}}}$   
 with  $\varepsilon_n \sim \text{Student-}t$

The Student-t parameters are hard to guess

- mean (maybe 0)
- variance (maybe 1)
- degrees of freedom (maybe 2)

even harder to guess!

NB: Important to note that outliers are present and that a Gaussian noise model may not be the most appropriate choice here.

7.

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

$$b) \quad \left. \frac{d\mathcal{L}(a)}{da} \right|_{a_{ML}} = \sum_{n=1}^N x_n (y_n - a_{ML} x_n) = 0 \iff a_{ML} = \frac{\sum_{n=1}^N x_n y_n}{\sum_{n=1}^N x_n^2}$$

8.

$$a) \quad \mathcal{L}(a) = \log P(\{y_n\}_{n=1}^N, \{z_n\}_{n=1}^N | \{x_n\}_{n=1}^N, a) = \sum_{n=1}^N \log P(y_n, z_n | x_n, a)$$

where  $P(y_n, z_n | x_n, a) = \mathcal{N}\left(\begin{bmatrix} y_n \\ z_n \end{bmatrix} ; \begin{bmatrix} ax_n \\ x_n \end{bmatrix}, \Sigma = \begin{pmatrix} 1 & 1/2 \\ 1/2 & 1 \end{pmatrix}^{-1}\right)$

$$\begin{aligned}
 \Rightarrow \mathcal{L}(a) &= \sum_{n=1}^N -\frac{1}{2} \left( \begin{bmatrix} y_n \\ z_n \end{bmatrix} - \begin{bmatrix} ax_n \\ x_n \end{bmatrix} \right)^T \begin{pmatrix} 1 & 1/2 \\ 1/2 & 1 \end{pmatrix} \left( \begin{bmatrix} y_n \\ z_n \end{bmatrix} - \begin{bmatrix} ax_n \\ x_n \end{bmatrix} \right) - \frac{N}{2} \log(\det(2\pi\Sigma)) \\
 &= -\frac{1}{2} \left( \sum_n (y_n - ax_n)^2 + \sum_n (y_n - ax_n)(z_n - x_n) + \sum_n (z_n - x_n)^2 + N \log(\det(2\pi\Sigma)) \right)
 \end{aligned}$$

$$\begin{aligned}
 b) \quad & \frac{d\mathcal{L}(a)}{da} = \sum_n \overbrace{(y_n - ax_n)x_n}^{a \text{ bit of info from just observing } y_i} + \frac{1}{2} \sum_n \overbrace{(z_n - x_n)x_n}^{\text{extra bit from observing } z \text{'s}} = 0 \iff \\
 & \iff \sum_n y_n x_n + \frac{1}{2} \sum_n z_n x_n - \frac{1}{2} \sum_n x_n^2 - a \sum_n x_n^2 = 0 \iff \\
 & \iff a = \left( \sum_n y_n x_n + \underbrace{\frac{1}{2} \sum_n (z_n - x_n)x_n}_{\substack{\text{new contribution from observing } z \text{'s}}} \right) / \sum_n x_n^2 \quad \begin{matrix} (\text{max likelihood}) \\ \text{estimate} \end{matrix}
 \end{aligned}$$

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  $y_n$ , and so observing  $z_n$  reveals information about the noise  $\epsilon_n$  and allows more accurate identification of  $a$ .

9. Info:  $P(y|m, c, x) = N(y; mx + c, 1)$   $P(m) = N(m; 0, 1)$   $P(c) = \text{u}N(c; 0, 1)$

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

$$\propto P(y|m, c, x) \cdot P(m) \cdot P(c) = \text{product of Gaussian functions in } m \text{ & } c$$

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

Tactic:

1) Substitute in for these densities.

2) Identify mean and covariance by comparing coefficients.

$$1) P(m, c | y, x) \propto \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y-mx-c)^2} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}m^2} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}c^2} \propto$$

$$\propto \exp\left(-\frac{1}{2}y^2 - \frac{1}{2}m^2(x^2+1) - \frac{1}{2}c^2 \times 2 - mcx + \cancel{y(mx+c)}\right) \propto$$

$$\propto \exp\left(-\frac{1}{2}\left(\begin{bmatrix} m \\ c \end{bmatrix} - \mu^{\text{POST}}\right)\left(\Sigma^{\text{POST}}\right)^{-1}\left(\begin{bmatrix} m \\ c \end{bmatrix} - \mu^{\text{POST}}\right)\right) \propto$$

$$\propto \exp\left(-\frac{1}{2}\left(\begin{bmatrix} m \\ c \end{bmatrix}^T \left(\Sigma^{\text{POST}}\right)^{-1} \begin{bmatrix} m \\ c \end{bmatrix} + \cancel{\begin{bmatrix} m \\ c \end{bmatrix}^T \left(\Sigma^{\text{POST}}\right)^{-1} \mu^{\text{POST}}}\right)\right)$$

← inverse covariance times mean  
 $\begin{bmatrix} m \\ c \end{bmatrix}^T \begin{bmatrix} y \\ x \end{bmatrix}$

$$\Rightarrow \Sigma^{\text{POST}} = \begin{bmatrix} x^2+1 & x \\ x & 2 \end{bmatrix}^{-1} = \frac{1}{2(x^2+1) - x^2} \begin{pmatrix} 2 & -x \\ -x & x^2+1 \end{pmatrix} = \frac{1}{x^2+2} \begin{pmatrix} 2 & -x \\ -x & x^2+1 \end{pmatrix}$$

$$\Rightarrow \text{Since } (\Sigma^{\text{POST}})^{-1} \mu^{\text{POST}} = \begin{bmatrix} y \\ x \end{bmatrix} \Rightarrow \mu^{\text{POST}} = \frac{1}{x^2+2} \begin{pmatrix} 2 & -x \\ -x & x^2+1 \end{pmatrix} \begin{bmatrix} y \\ x \end{bmatrix} =$$

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

b) Check these expressions match with those in lectures:

$$\underline{y} = \tilde{\underline{x}} \underline{w} + \underline{\varepsilon}, \text{ where } \underline{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \underline{\underline{I}}), \underline{w} \sim \mathcal{N}(0, \lambda^{-1} \underline{\underline{I}})$$

$$P(\underline{w} | \text{Data}, \sigma^2, \lambda) = \mathcal{N}(\underline{w}; \underline{\mu}_{\text{wid}}, \underline{\Sigma}_{\text{wid}})$$

$$\text{where } \underline{\Sigma}_{\text{wid}} = (\lambda \underline{\underline{I}} + \frac{1}{\sigma^2} \tilde{\underline{x}}^T \tilde{\underline{x}})^{-1}$$

$$\underline{\mu}_{\text{wid}} = \sum_{\text{wid}} \frac{1}{\sigma^2} \tilde{\underline{x}}^T \underline{y}$$

In this exercise  $\sigma^2 = \lambda = 1$  (prior and observation noise are unit variance)

$$\begin{aligned} \tilde{\underline{x}} &= [x, 1] \\ \underline{w} &= [m] \end{aligned} \Rightarrow \begin{aligned} \underline{\Sigma}_{\text{wid}} &= (\underline{\underline{I}} + [\underline{x}] [\underline{x}, 1]^T)^{-1} = \begin{pmatrix} x^2 + 1 & x \\ x & 2 \end{pmatrix}^{-1} \checkmark \\ \underline{\mu}_{\text{wid}} &= \begin{pmatrix} x^2 + 1 & x \\ x & 2 \end{pmatrix}^{-1} [\underline{x}] \underline{y} \checkmark \end{aligned}$$

c) We just need to plug this settings into  $\underline{\Sigma}^{\text{POST}}$  and  $\underline{\mu}^{\text{POST}}$  calculated in a).

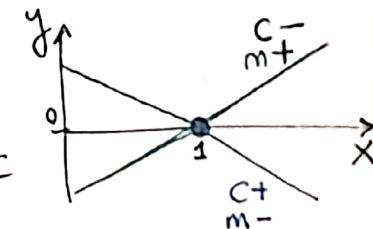
$$\text{i)} \begin{cases} x=0 \\ y=0 \end{cases} \Rightarrow \underline{\Sigma}^{\text{POST}} = \begin{bmatrix} 1 & 0 \\ 0 & 1/2 \end{bmatrix}, \quad \underline{\mu}^{\text{POST}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

data doesn't tell us anything about  $m$ :  $y = mx + c$  and  $x=0$ . This is why posterior mean and variance of  $m$  stays at prior (mean 0, variance 1). Data tells us that  $c$  is likely to be small ( $\approx 0$ ), but the observation noise is quite large so only weak evidence for this. This is why posterior mean is 0 and posterior variance is  $1/2$  (prior variance is 1).

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

The data can be explained by either positive  $c$  and negative  $m$  or negative  $c$  & positive  $m$ : hence posterior covariance between  $m$  &  $c$  is negative ( $-1/3$ ).

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



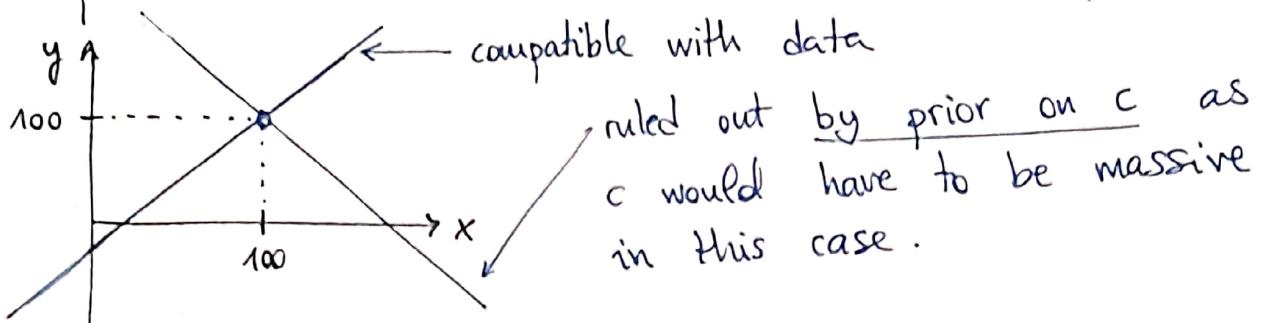
$$\text{iii) } \begin{aligned} x &= 100 \\ y &= 100 \end{aligned} \Rightarrow \sum^{\text{post}} = \begin{bmatrix} \frac{2}{100^2} & -\frac{1}{100} \\ -\frac{1}{100} & 1 \end{bmatrix}, \quad M^{\text{post}} = \begin{bmatrix} 1 \\ \frac{1}{100} \end{bmatrix}$$

The data  $y=100, x=100$  rules out negative  $m$ :

We are therefore quite certain that  $m$  is close to 1

(posterior mean very close to this value and variance  $\approx \frac{2}{100^2}$ )

The data give almost no information about  $c$  though with posterior mean and variance close to the prior.



c.2) Predictive mean: use  $y^* = mx^* + c + \varepsilon^*$   
 now average w.r.t.  $P(m, c, \varepsilon^* | x, y)$ :

$$\mathbb{E}_{P(m, c, \varepsilon^* | x, y)} [mx^* + c + \varepsilon^*] = \mu_m^{\text{post}} x^* + \mu_c^{\text{post}} + 0$$

$\uparrow$  independent from  $x$  and  $y$

i) predictive mean = 0 (horizontal line)

ii) predictive mean = 0 (horizontal line)

iii) predictive mean =  $\frac{1}{100^2+2} (100^2 x^* + 100)$

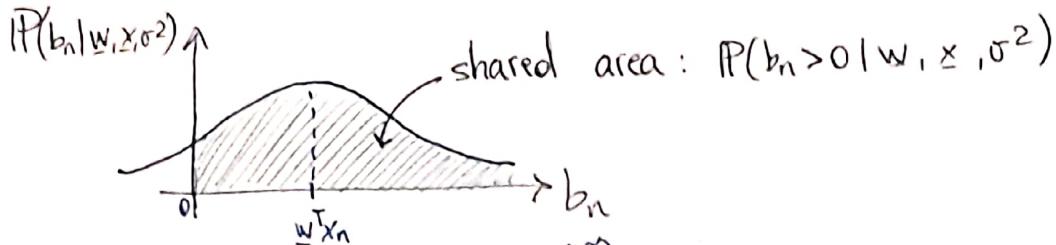
line with gradient very close to 1  
 and a small positive intercept.

10.

a) Info:  $y_n = H(\underline{w}^T \underline{x}_n + \varepsilon_n)$  where  $\varepsilon_n \sim \mathcal{N}(0, \sigma^2)$

Let  $b_n = \underline{w}^T \underline{x}_n + \varepsilon_n$ ,  $b_n \sim \mathcal{N}(\underline{w}^T \underline{x}_n, \sigma^2)$

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



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

Let  $v_n = -\frac{b_n - \underline{w}^T \underline{x}_n}{\sigma^2}$ , and transform to integral over this quantity ( $v_n$  will be distributed to a standard normal)

- when  $b_n = 0 \Rightarrow v_n = \frac{\underline{w}^T \underline{x}_n}{\sigma^2}$

- when  $b_n = \infty \Rightarrow v_n = -\infty$

$$\Rightarrow P(y_n = 1 | \underline{w}, \underline{x}_n, \sigma^2) = \int_{-\infty}^{\frac{\underline{w}^T \underline{x}_n}{\sigma^2}} \mathcal{N}(v_n; 0, 1) dv_n = \text{CDF}\left(\frac{\underline{w}^T \underline{x}_n}{\sigma^2}\right)$$

b) If  $\sigma^2 \rightarrow \infty \Rightarrow P(y_n = 1 | \underline{x}_n, \underline{w}, \sigma^2) = \text{CDF}(0) = 1/2$

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

11.

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

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

$$\hat{\underline{w}}_{\max} = \operatorname{argmax}_{\hat{\underline{w}} \in \hat{\underline{w}}_1, \dots, \hat{\underline{w}}_K} \hat{\underline{w}}^\top \underline{x}$$

$i \in \{1, \dots, K\}$

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

$$\Rightarrow \text{As } \beta \rightarrow \infty, e^{\beta (\hat{\underline{w}}_i - \hat{\underline{w}}_{\max})} \rightarrow \begin{cases} 0 & \text{if } \hat{\underline{w}}_i \neq \hat{\underline{w}}_{\max} \\ \infty & \text{if } \hat{\underline{w}}_i = \hat{\underline{w}}_{\max} \end{cases}$$

*shouldn't it be 1?*

$$\Rightarrow y_i(\underline{x}, \underline{w}) \rightarrow \begin{cases} 0 & \text{if } i \neq i_{\max} \\ 1 & \text{if } i = i_{\max} \end{cases}$$

This is a zero/one encoding of the argmax function, hence the general name "softmax":

i.e. as  $\beta \rightarrow \infty$   $y_i(\underline{x}, \underline{w}) \rightarrow \mathbb{I}(\operatorname{argmax}_K \hat{\underline{w}}_K^\top \underline{x} = i)$

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

b)  $y_i = P(t_i = 1 | \underline{w}, \underline{x})$

$$P(\{t^{(n)}\}_{n=1}^N | \underline{w}, \{x^{(n)}\}_{n=1}^N) = \prod_{n=1}^N P(t^{(n)} | \underline{w}, x^{(n)}) =$$

$$= \prod_{n=1}^N \prod_{k=1}^K y_k(x^{(n)}, \underline{w})^{t_k^{(n)}} = \exp\left(\sum_n \sum_k t_k^{(n)} \log y_k(x^{(n)}, \underline{w})\right)$$

$\Rightarrow$  cost function could be  $\operatorname{argmax}_{\underline{w}} \sum_{n=1}^N \sum_{k=1}^K t_k^{(n)} \log y_k(x^{(n)}, \underline{w})$

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

$$y_1(\underline{x}, \underline{w}) = \frac{e^{\underline{w}_1^\top \underline{x}}}{e^{\underline{w}_1^\top \underline{x}} + e^{\underline{w}_2^\top \underline{x}}} = \frac{1}{1 + e^{(\underline{w}_2 - \underline{w}_1)^\top \underline{x}}} = \frac{1}{1 + e^{\underline{v}^\top \underline{x}}}$$

where  $\underline{v} = \underline{w}_2 - \underline{w}_1$ . This is equivalent to the expression for logistic regression applied to a two-class classification problem. However, notice that although the softmax and logistic functions are identical for  $K=2$  classes, the parameterisation is different. The softmax version is over-parameterised having two sets of parameters whose difference affects the input-output function. For this reason the parameters of the softmax are not identifiable: adding the same vector to each weight  $w_k \leftarrow w_k + b$  causes no change in the input-output function.