

## Regression

## Regresssion

- A regression problem is one that identifies patterns and regularities in the mapping b/w ip variables  $\underline{x}^{(n)}$  and a corresponding continuous op variable  $y^{(n)}$  from a training data set  $D = \{(\underline{x}^{(n)}, y^{(n)})\}_{n=1}^N$
- (Note  $\underline{x}^{(n)} = [x_1^{(n)}, \dots, x_d^{(n)}]^T \in \mathbb{R}^d$  and  $y_n \in \mathbb{R}$ )
- A regression method typically returns a conditional distribution  $p(y| \underline{x})$  for the unknown  $y$  associated w/ a new test ip  $\underline{x}_t \rightarrow$  useful for predictions / interpretations (find correlation)
- ↳ The characteristics of  $p(y| \underline{x})$  should capture patterns observed in the data (e.g. changes in mean/variance)
- We use a model to encode, as function of some parameters  $\underline{\theta}$ , the type of patterns we expect to find in the data.

## Linear regression model

- For the linear regression model, we assume that the relationship b/w ip  $\underline{x}^{(n)}$  and op  $y^{(n)}$  in  $p(y| \underline{x})$  is linear, w/  $\underline{\theta}$  as parameters specifying this conditional distribution.
- For the 1D case, we have  $\underline{\theta} = \{\sigma^2, w_0, w_1\}$

$$y^{(n)} = w_0 + w_1 x^{(n)} + \varepsilon^{(n)}, \quad \varepsilon^{(n)} \sim N(0, \sigma^2)$$

$y^{(n)}$  is the sum of a const. and a Gaussian RV  $\rightarrow y^{(n)}$  is Gaussian distributed.

$$\begin{aligned} E[y^{(n)}] &= w_0 + w_1 x^{(n)} + E[\varepsilon^{(n)}] = w_0 + w_1 x^{(n)} & \text{Var}[y^{(n)}] &= E[(y^{(n)} - E[y^{(n)}])^2] = E[\varepsilon^{(n)2}] = \sigma^2 \\ \therefore p(y^{(n)} | \underline{x}^{(n)}, \underline{\theta}) &= N(y^{(n)} | w_0 + w_1 x^{(n)}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2} \frac{(y^{(n)} - w_0 - w_1 x^{(n)})^2}{\sigma^2}\right] \end{aligned}$$

- For higher dimensions, we have  $\underline{\theta} = \{\sigma^2, \underline{w} = [w_0 \dots w_d]^T\}$

$$y^{(n)} = w_0 + w_1 x_1^{(n)} + \dots + w_d x_d^{(n)} = [\underline{w} \dots w_d] \begin{bmatrix} x_1^{(n)} \\ x_2^{(n)} \\ \vdots \\ x_d^{(n)} \end{bmatrix} + \varepsilon_n = \underline{w}^T \underline{x}^{(n)} + \varepsilon_n, \quad \varepsilon_n \sim N(0, \sigma^2)$$

where  $\underline{x} = [1, \underline{x}^{(n)}]^T$  is the augmented ip

$y^{(n)}$  is the sum of a const. and a Gaussian RV  $\rightarrow y^{(n)}$  is Gaussian distributed.

$$\begin{aligned} E[y^{(n)}] &= \underline{w}^T \underline{x}^{(n)} + E[\varepsilon^{(n)}] = \underline{w}^T \underline{x}^{(n)} & \text{Var}[y^{(n)}] &= E[(y^{(n)} - E[y^{(n)}])^2] = E[\varepsilon^{(n)2}] = \sigma^2 \\ \therefore p(y^{(n)} | \underline{x}^{(n)}, \underline{\theta}) &= N(y^{(n)} | \underline{w}^T \underline{x}^{(n)}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2} \frac{(y^{(n)} - \underline{w}^T \underline{x}^{(n)})^2}{\sigma^2}\right] \end{aligned}$$

## Maximum likelihood estimate (MLE).

- We find the MLE for  $\underline{\theta}$  by max. the likelihood function  $p(y^{(1)} \dots y^{(N)} | \underline{x}^{(1)} \dots \underline{x}^{(N)}, \underline{\theta})$  wrt  $\underline{\theta}$ .

$$\underline{\theta}^{ML} = \underset{\underline{\theta}}{\operatorname{arg\,max}} p(y^{(1)} \dots y^{(N)} | \underline{x}^{(1)} \dots \underline{x}^{(N)}, \underline{\theta})$$

- In practice, we use the log-likelihood function  $L(\underline{\theta}) = \log p(y^{(1)} \dots y^{(N)} | \underline{x}^{(1)} \dots \underline{x}^{(N)}, \underline{\theta})$

↳ we get the same soln since logarithm is monotonically increasing

$$\underline{\theta}^{ML} = \underset{\underline{\theta}}{\operatorname{arg\,max}} p(y^{(1)} \dots y^{(N)} | \underline{x}^{(1)} \dots \underline{x}^{(N)}, \underline{\theta}) = \underset{\underline{\theta}}{\operatorname{arg\,max}} L(\underline{\theta})$$

↳ taking the log can prevent overflow/underflow problems when computing products.

# For Personal Use Only -bkwk2

- Assuming that each data pt is independent, then

$$p(\underline{y} | \underline{x}, \theta) = p(y^1 | x^{(1)}, \dots, y^{(N)} | x^{(N)}, \theta) = \prod_{n=1}^N p(y^{(n)} | x^{(n)}, \theta) = \prod_{n=1}^N N(y^{(n)} | \underline{w}^T \underline{x}^{(n)}, \sigma^2)$$

The loglikelihood  $L(\theta)$  is thus given by

$$L(\theta) = \log p(y^1 | x^{(1)}, \dots, y^{(N)} | x^{(N)}, \theta) = \sum_{n=1}^N \log N(y^{(n)} | \underline{w}^T \underline{x}^{(n)}, \sigma^2)$$

$$= \sum_{n=1}^N \left[ -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (y^{(n)} - \underline{w}^T \underline{x}^{(n)})^2 \right] = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\underline{y} - \underline{\tilde{x}}\underline{w})^T (\underline{y} - \underline{\tilde{x}}\underline{w})$$

$$\boxed{L(\theta) = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\underline{w}^T \underline{\tilde{x}}\underline{\tilde{x}}\underline{w} - 2\underline{y}^T \underline{\tilde{x}}\underline{w} + \underline{y}^T \underline{y})}$$

where  $\underline{y} = [y_1 \dots y_N]^T$  and  $\underline{\tilde{x}} = \begin{bmatrix} \underline{x}^{(1)} \\ \vdots \\ \underline{x}^{(N)} \end{bmatrix}$

- Differentiating  $L(\theta)$  wrt  $\underline{w}$ ,

$$\frac{\partial L(\theta)}{\partial \underline{w}} = \frac{\partial}{\partial \underline{w}} \left( -\frac{1}{2\sigma^2} (\underline{w}^T \underline{\tilde{x}}\underline{\tilde{x}}\underline{w} - 2\underline{y}^T \underline{\tilde{x}}\underline{w}) \right) = -\frac{\underline{\tilde{x}}\underline{\tilde{x}}\underline{w}}{\sigma^2} + \frac{\underline{\tilde{x}}\underline{y}}{\sigma^2} = 0$$

$$\therefore \boxed{\underline{w} = (\underline{\tilde{x}}^T \underline{\tilde{x}})^{-1} \underline{\tilde{x}}^T \underline{y}}$$

Differentiating  $L(\theta)$  wrt  $\sigma^2$ ,

$$\frac{\partial L(\theta)}{\partial \sigma^2} = \frac{\partial}{\partial \sigma^2} \left( -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\underline{y} - \underline{\tilde{x}}\underline{w})^T (\underline{y} - \underline{\tilde{x}}\underline{w}) \right) = -\frac{N}{2\sigma^2} + \frac{1}{2\sigma^4} (\underline{y} - \underline{\tilde{x}}\underline{w})^T (\underline{y} - \underline{\tilde{x}}\underline{w}) = 0$$

$$\therefore \boxed{\sigma^2 = \frac{1}{N} (\underline{y} - \underline{\tilde{x}}\underline{w})^T (\underline{y} - \underline{\tilde{x}}\underline{w})}$$

+ when  $N \leq D$ , the MLE  $\underline{w} = (\underline{\tilde{x}}^T \underline{\tilde{x}})^{-1} \underline{\tilde{x}}^T \underline{y}$  is not defined, as  $\underline{\tilde{x}}^T \underline{\tilde{x}}$  is not invertible

→ we have many values of  $\underline{w}$  that fit the training data equally well w/ zero error.

## Non-linear (basis function) regression

- Linear regression can model non-linear relationships by replacing  $\underline{x}^{(n)}$  w/ some non-linear function of the i/p's  $\Phi(\underline{x}) = [\Phi_1(\underline{x}) \dots \Phi_M(\underline{x})]^T$  (so  $\underline{w}^T \underline{x}^{(n)} \rightarrow \underline{w}^T \Phi(\underline{x}^{(n)})$ )

- Inference does not change, we just replace each  $\underline{x}^{(n)}$  w/ the new  $\Phi(\underline{x}^{(n)})$

- Common basis functions  $\Phi_m(\underline{x})$  include:

↳ polynomial for 1D data:  $\Phi_m(\underline{x}^{(n)}) = \underline{x}^{(n)m}$

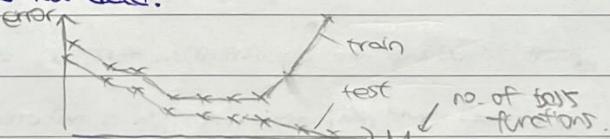
↳ Gaussian radial basis functions:  $\Phi_m(\underline{x}^{(n)}) = \exp \left[ -\frac{(\underline{x}^{(n)} - \underline{c}_m)^T (\underline{x}^{(n)} - \underline{c}_m)}{2s^2} \right]$

↳ sigmoidal basis functions:  $\Phi_m(\underline{x}^{(n)}) = \sigma \left( \frac{(\underline{x}^{(n)} - \underline{c}_m)^T (\underline{x}^{(n)} - \underline{c}_m)}{s^2} \right)$ ,  $\sigma(x) = \frac{1}{1 + \exp(-x)}$

↳ Linear in terms of transformed i/p's  $\Phi(\underline{x}^{(n)})$  but non-linear in terms of original i/p's  $\underline{x}^{(n)}$ .

overfitting.

- A large no. of basis functions can lead to overfitting – the model fits the training data well but it performs poorly on new test data.



- As  $M$  increases, the magnitude of the entries of the MLE of  $\underline{w}$  become large → large oscillations in the model's predictions → solve this by using a prior to favor the entries of  $\underline{w}$  to be small.

# For Personal Use Only -bkwk2

## Bayesian linear regression

### Multivariate Gaussian distribution

- The density of a D-dimensional vector  $\underline{x} \in \mathbb{R}^D$  is

$$N(\underline{x} | \underline{m}, \underline{\Sigma}) = \frac{1}{(2\pi)^{D/2} |\underline{\Sigma}|} \exp\left[-\frac{1}{2} (\underline{x}-\underline{m})^T \underline{\Sigma}^{-1} (\underline{x}-\underline{m})\right] \propto \exp\left[-\frac{1}{2} \underline{x}^T \underline{\Sigma}^{-1} \underline{x} + \underline{m}^T \underline{\Sigma}^{-1} \underline{x}\right]$$

i.e. the density is prop. to the exponential of a quadratic function of  $\underline{x}$ , w/ normalization const.  $Z$

$$Z = \int (2\pi)^D |\underline{\Sigma}| \exp\left(-\frac{1}{2} \underline{m}^T \underline{\Sigma}^{-1} \underline{m}\right)$$

where  $\underline{m} \in \mathbb{R}^D$  is the mean vector and  $\underline{\Sigma} \in \mathbb{R}^{D \times D}$  is the covariance matrix

- The general Gaussian  $p(\underline{x}) = N(\underline{x} | \underline{m}, \underline{\Sigma})$  can be obtained from the standard Gaussian  $p(\underline{z} | \underline{0}, I)$

$$\underline{x} = \underline{m} + \underline{\Sigma}^{1/2} \underline{z}$$

where  $\underline{\Sigma}^{1/2}$  is the matrix square root, defined as  $\underline{\Sigma} = \underline{\Sigma}^{1/2} (\underline{\Sigma}^{1/2})^T$  ↗ Cholesky decomposition

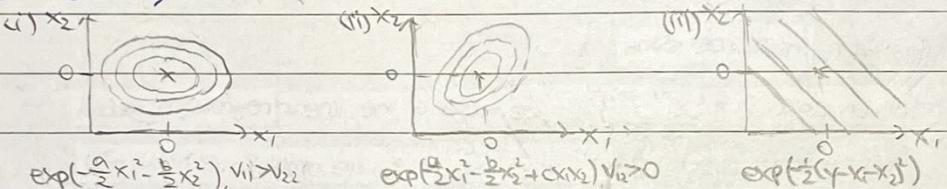
### Effect of parameters on the density function

- For the general Gaussian density,

$$p(\underline{x}) = N(\underline{x} | \underline{m}, \underline{\Sigma}) = \frac{1}{(2\pi)^{D/2} |\underline{\Sigma}|} \exp\left[-\frac{1}{2} (\underline{x}-\underline{m})^T \underline{\Sigma}^{-1} (\underline{x}-\underline{m})\right]$$

the parameter  $\underline{m}$  determines the mode (location) and  $\underline{\Sigma}$  scales+rotates the space.

- For the  $D=2$  case,



$$(i) \exp\left(-\frac{a}{2}x_1^2 - \frac{b}{2}x_2^2\right), V_{11} > V_{22}$$

$$(ii) \exp\left(-\frac{a}{2}x_1^2 - \frac{b}{2}x_2^2 + cx_1x_2\right), V_{11} > 0$$

$$(iii) \exp\left(-\frac{1}{2}(y-x_1-x_2)^2\right) = c \rightarrow x_1 + x_2 = y + k$$

### Linear combination of Gaussian RVs

- Let  $p(\underline{x}) = N(\underline{x} | \underline{0}, \underline{\Sigma})$  and  $p(\underline{e}) = N(\underline{e} | \underline{0}, \underline{\Sigma}_e)$ . Consider the transformation

$$\underline{y} = \underline{W} \underline{x} + \underline{e}$$

As linear combinations of independent Gaussian RVs are Gaussian,  $p(\underline{y})$  is Gaussian,  $p(\underline{y}) = N(\underline{y} | \underline{m}_y, \underline{\Sigma}_y)$

$$\hookrightarrow \underline{m}_y = E[\underline{y}] = \underline{W} E[\underline{x}] + E[\underline{e}] = \underline{0} \rightarrow \underline{m}_y = \underline{0}$$

$$\hookrightarrow \underline{\Sigma}_y = E[\underline{y}\underline{y}^T] - E[\underline{y}]E[\underline{y}]^T = E[(\underline{W}\underline{x} + \underline{e})(\underline{W}\underline{x} + \underline{e})^T] = E[\underline{W}\underline{x}\underline{x}^T \underline{W}^T + \underline{e}\underline{e}^T + \underline{W}\underline{x}\underline{e}^T + \underline{e}\underline{x}^T]$$

$$= \underline{W} E[\underline{x}\underline{x}^T] \underline{W}^T + E[\underline{e}\underline{e}^T] E[\underline{x}\underline{x}^T] + E[\underline{x}\underline{x}^T] E[\underline{e}\underline{e}^T] + E[\underline{e}\underline{e}^T] = \underline{W} \underline{V}_1 \underline{W}^T + \underline{V}_2 \rightarrow \underline{V}_y = \underline{W} \underline{V}_1 \underline{W}^T + \underline{V}_2$$

\* Note that  $\underline{x}$  and  $\underline{e}$  are assumed to be independent, so  $E[\underline{x}\underline{e}^T] = E[\underline{x}]E[\underline{e}^T]$  etc.

# For Personal Use Only -bkwk2

Product of Gaussian densities

- Let  $p(x) = N(x | m_1, \Sigma_1)$  and  $q(x) = N(x | m_2, \Sigma_2)$ , consider the product

$$t(x) \propto p(x) q(x)$$

since the product of exponentials of quadratic functions is also the exponential of a quadratic function,  $t(x)$  is also Gaussian,  $t(x) = N(x | m_3, \Sigma_3)$

$$\begin{aligned} t(x) \propto p(x) q(x) &= N(x | m_1, \Sigma_1) N(x | m_2, \Sigma_2) = \frac{1}{\sqrt{(2\pi)^n |\Sigma_1|}} \frac{1}{\sqrt{(2\pi)^n |\Sigma_2|}} \exp\left[-\frac{1}{2}(x-m_1)^T \Sigma_1^{-1} (x-m_1) - \frac{1}{2}(x-m_2)^T \Sigma_2^{-1} (x-m_2)\right] \\ &\propto \exp\left[-\frac{1}{2} x^T (\Sigma_1^{-1} + \Sigma_2^{-1}) x + (m_1^T \Sigma_1^{-1} + m_2^T \Sigma_2^{-1}) x\right] \end{aligned}$$

Noting that  $t(x) = N(x | m_3, \Sigma_3) \propto \exp(-\frac{1}{2} x^T \Sigma_3^{-1} x + m_3^T \Sigma_3^{-1} x)$ ,

$$\begin{aligned} \Sigma_1^{-1}, \Sigma_2^{-1} \text{ are} \\ \text{"inverted parameters"} \\ \rightarrow \text{empty sum} \quad \boxed{\Sigma_3^{-1} = (\Sigma_1^{-1} + \Sigma_2^{-1})} &\rightarrow \Sigma_3 = (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1} \\ \boxed{\Sigma_3^{-1} = \Sigma_1^{-1} + \Sigma_2^{-1}} &\rightarrow m_3 = \Sigma_3 (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1} \end{aligned}$$

completing the square.

- Let  $p(x) = N(x | m, \Sigma) \propto \exp\left[-\frac{1}{2} x^T \Sigma^{-1} x + m^T \Sigma^{-1} x\right]$  and we are given an exponential of a quadratic function  $q(x) \propto \exp\left[-\frac{1}{2} x^T B x + b^T x\right]$  to match w/.

$$\boxed{B = \Sigma^{-1}} \quad \boxed{\Sigma = B^{-1}} \quad \boxed{b^T = m^T \Sigma^{-1}} \quad \rightarrow \quad \boxed{m = \Sigma b}$$

Bayesian inference

- Using a large no. of basis functions can lead to overfitting - the model fits the training data well but performs poorly on new test data  $\rightarrow \|y^{(M)}\|$  becomes large  $\rightarrow$  large oscillations in model predictions.
- We can favour smooth solns using priors to enforce  $\theta$  to be small.

- The posterior distribution for  $\theta$ ,  $p(\theta | y, x)$  is obtained by Bayes rule

$$p(\theta | y, x) = \frac{p(y | x, \theta) p(\theta | x)}{p(y | x)}$$

where  $p(y | x, \theta) = \prod_{i=1}^n p(y^{(i)} | x^{(i)}, \theta)$  is the likelihood,  $p(\theta | x)$  is the model evidence

and  $p(\theta | x) = p(\theta)$  is the prior distribution.

- The predictive distribution for  $y_*$  given a new corresponding  $x_*$  is

$$p(y_* | x_*, y, x) = \int p(y_* | x_*, y, x, \theta) d\theta = \int p(y_* | x_*, \theta) p(\theta | y, x) d\theta$$

i.e. weighted average of the predictions  $p(y_* | x_*, \theta)$ , over the posterior distribution  $p(\theta | y, x)$ .

- Exact inference is possible if prior and noise distributions are Gaussian.

# For Personal Use Only -bkwk2

## Bayesian linear regression

- Given data  $D = \{(x^{(n)}, y^{(n)})\}_{n=1}^N$ , we assume the linear regression model

$$y^{(n)} = w^T x^{(n)} + \epsilon_n, \quad \epsilon_n \sim N(0, \sigma^2)$$

To simplify inference, assume  $\sigma^2$  is known — the only unknown is  $w$ .

- Choose the prior for  $w$  to be a zero-mean isotropic Gaussian precision matrix

$$p(w) = N(w|0, \lambda I) \propto \exp\left[-\frac{1}{2} w^T \lambda I w\right]$$

Recall that the likelihood under Gaussian noise is

$$p(y|\Sigma, w) = \prod_{n=1}^N N(y^{(n)} | \Sigma^{-1} w, \sigma^2) \propto \exp\left[-\frac{\Sigma^{-1} w^T \Sigma w}{2\sigma^2} + \frac{y^T \Sigma w}{\sigma^2}\right]$$

The posterior density is the product of two Gaussian densities  $\rightarrow$  posterior is Gaussian  $N(m, V)$

$$p(w|y, \Sigma, \sigma^2) \propto p(y|\Sigma, w) p(w) \propto \exp\left[-\frac{1}{2} w^T \left(\frac{\Sigma^{-1} \Sigma}{\sigma^2} + \lambda I\right) w + \frac{y^T \Sigma w}{\sigma^2}\right]$$

$$\text{Therefore } V = \left(\frac{\Sigma^{-1} \Sigma}{\sigma^2} + \lambda I\right)^{-1}, \quad m = \frac{\Sigma^{-1} y}{\sigma^2}$$

$$(\text{Alternatively, note } V = \frac{\Sigma^{-1} \Sigma}{\sigma^2} + (\lambda I)^{-1}, \quad m^T V = \frac{y^T \Sigma}{\sigma^2} + \sigma^2 (\lambda I)^{-1})$$

- As  $\lambda \rightarrow 0$  (prior is flat),  $V^{MAP} = m = w^{ML}$  i.e. ML and MAP solutions are identical.

## Bayesian predictive distribution

- The predictive distribution for  $y_*$  given a new corresponding  $\tilde{x}_*$  is

$$p(y_*|\tilde{x}_*, y, \Sigma) = \int p(y_*|w, \tilde{x}_*) p(w|y, \Sigma) dw = \int N(y_*|w^T \Sigma^{-1} \tilde{x}_*, \sigma^2) N(w|m, V) dw$$

- Instead of directly computing the integral, we have  $y_* \sim \Sigma^{-1} \tilde{x}_* + e^*$ , where  $e^* \sim N(0, V)$ ,

$$p(y_*|\tilde{x}_*, y, \Sigma) = N(y_*|m_*, V_*)$$

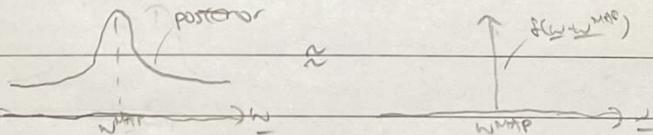
$$\text{where } m_* = E[y_*] = E[w^T \Sigma^{-1} \tilde{x}_*] + E[e^*] = m^T \tilde{x}_* \rightarrow m_* = m^T \tilde{x}_*$$

$$V_* = \text{Var}[y_*] = \text{Var}[w^T \Sigma^{-1} \tilde{x}_*] + \text{Var}[e^*] = \Sigma^{-1} \Sigma \tilde{x}_*^T + \sigma^2 \rightarrow V_* = \Sigma^{-1} \Sigma \tilde{x}_*^T + \sigma^2$$

- We can reduce overfitting and obtain confidence bands  $m_* \pm \frac{1}{2} V_*$  in our predictions.

## Maximum a priori (MAP) inference

- Assume that the posterior is well approximated by a pt mass at its mode



The predictive distribution becomes

$$p(y_*|\tilde{x}_*, y, \Sigma) = \int p(y_*|w, \tilde{x}_*) p(w|y, \Sigma) dw \approx \int p(y_*|w, \Sigma) \delta(w - w^{MAP}) dw = p(y_*|\tilde{x}_*, w^{MAP})$$

- MAP inference fails to generate confidence bands in the resulting prediction.

- The MAP soln  $w^{MAP}$  is obtained as follows

$$w^{MAP} = \arg \max_w p(w|y, \Sigma) = \arg \max_w p(y|w, \Sigma) p(w) = \arg \max_w [\log p(y|w, \Sigma) + \log p(w)]$$

$$\text{For } p(w) = N(w|0, \lambda I), \quad w^{MAP} = \arg \max_w [\log p(y|w, \Sigma) - \frac{1}{2} w^T \Sigma w] = (\Sigma^{-1} \Sigma + \lambda I)^{-1} \Sigma y$$

## Classification

### Classification

#### Classification

In regression, we are given a training set  $D = \{(\underline{x}^{(n)}, y^{(n)})\}_{n=1}^N$  consisting of i/p's  $\underline{x}^{(n)} = [x_0^{(n)} \dots x_D^{(n)}]^T$  and o/p's  $y^{(n)}$ , but the o/p's are now discrete:  $y^{(n)} \in \{1, 2, \dots, C\}$  where  $C$  is the no. of classes/categories.

- In classification, we aim to identify a partition of the i/p space into  $C$  decision regions, one for each class. Each new i/p is assigned the class of its corresponding decision region, along w/ a measure of confidence.
- The decision regions are separated by decision boundaries — these are pts of which two classes have equal predictive probability.

#### Problems with methods for regression

- A way of using methods for regression for classification is to let the o/p  $\underline{y}^{(n)}$  be a  $C$ -dimensional vector w/ onehot encoding of the class for  $\underline{x}^{(n)}$  (if the class is  $c$ ,  $y_i^{(n)} = 1\{i=c\}$ ).
- We can then solve  $C$  linear regression problems, one for each class,

$$\underline{w} = (\underline{\underline{x}}^T \underline{\underline{x}})^{-1} \underline{\underline{x}}^T \underline{\underline{y}}$$

where  $\underline{\underline{y}} = \begin{bmatrix} \underline{y}_1^{(1)} \\ \vdots \\ \underline{y}_C^{(1)} \end{bmatrix}$  and we predict the class w/ the highest entry of  $\underline{\underline{x}}^T \underline{w}$ .

- For the binary case, we could set a simple threshold to determine the class (positive class label if the o/p of the regression model is above the threshold and vice versa).
- These methods have issues of underrepresenting certain classes or varying the threshold (underlying hypothesis) as we collect more data pts.

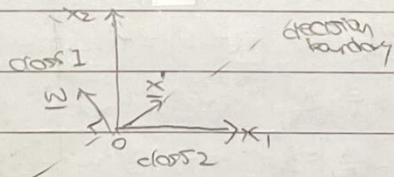
#### Deterministic linear classification

- Deterministic linear classification works by mapping the o/p of the linear model into discrete class labels. For the binary case, assume  $y^{(n)} \in \{0, 1\}$ , then we can define

$$y^{(n)} = H(\underline{w}^T \underline{x}^{(n)})$$

where  $H(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{otherwise} \end{cases}$  is the heaviside step function.

- The decision boundary is where  $\underline{w}^T \underline{x}^{(n)} = 0 \rightarrow \underline{w}$  orthogonal to  $\underline{x}^{(n)}$



- This model gives deterministic (hard) predictions — misclassification errors are not allowed and inference is hard (cannot use gradient-based methods to find the MLE)

↑  
step function is  
not differentiable

# For Personal Use Only -bkwk2

## Probabilistic linear classification

- Probabilistic linear classification works by mapping the o/p of the linear model into class probabilities (between 0 and 1). For the binary case, we define

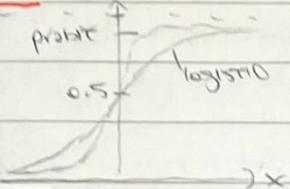
$$P(Y^{(n)} = 1 | \underline{x}, \underline{w}) = f(\underline{w}^T \underline{x})$$

where  $f(x)$  is a monotonically increasing function that maps  $\mathbb{R} \mapsto [0, 1]$ .

- Typical functions  $f(x)$  include:

↳ Logistic function:  $\sigma(x) = \frac{1}{1+e^{-x}}$

↳ probit function:  $L(\underline{x}) = \int_{-\infty}^{\underline{x}} N(z|0, 1) dz$



• Note that the functions are symmetric about  $(0, 0.5)$ , so  $f(-x) = 1 - f(x)$

- The likelihood function now allows for misclassification errors w/ some probability, and has a non-zero gradient that can be used to find the MLE.

## Logistic regression.

- In logistic regression, we use  $f(x) = \sigma(x) = \frac{1}{1+e^{-x}}$ .

(i)  $Y^{(n)} \in \{-1, 1\}$

$\{Y^{(n)} = 1\}$

$\{Y^{(n)} = 0\}$

$$P(Y^{(n)} | \underline{x}, \underline{w}) = \frac{1+Y^{(n)}}{2} \sigma(\underline{w}^T \underline{x}) + \frac{1-Y^{(n)}}{2} (1 - \sigma(\underline{w}^T \underline{x})) = \sigma(Y^{(n)} \underline{w}^T \underline{x})$$

$$L(\underline{w}) = \log P(Y | \underline{x}, \underline{w}) = \sum_{n=1}^N \log P(Y^{(n)} | \underline{x}, \underline{w}) = \sum_{n=1}^N \log \sigma(Y^{(n)} \underline{w}^T \underline{x})$$

Noting that  $\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x)) = \sigma(x)\sigma(-x)$ ,

$$\frac{dL(\underline{w})}{d\underline{w}} = \frac{d}{d\underline{w}} \sum_{n=1}^N \log \sigma(Y^{(n)} \underline{w}^T \underline{x}) = \sum_{n=1}^N \frac{d}{d\underline{w}} \log \sigma(Y^{(n)} \underline{w}^T \underline{x})$$

$$= \sum_{n=1}^N \frac{\partial}{\partial \underline{w}} \sigma(Y^{(n)} \underline{w}^T \underline{x}) \sigma(Y^{(n)} \underline{w}^T \underline{x}) (1 - \sigma(Y^{(n)} \underline{w}^T \underline{x})) \cdot Y^{(n)} \underline{x}$$

$$\frac{dL(\underline{w})}{d\underline{w}} = \sum_{n=1}^N Y^{(n)} (1 - \sigma(Y^{(n)} \underline{w}^T \underline{x})) \underline{x}$$

(ii)  $Y^{(n)} \in \{0, 1\}$

$$P(Y^{(n)} | \underline{x}, \underline{w}) = \sigma(\underline{w}^T \underline{x})^{Y^{(n)}} (1 - \sigma(\underline{w}^T \underline{x}))^{1-Y^{(n)}}$$

$$L(\underline{w}) = \log P(Y | \underline{x}, \underline{w}) = \sum_{n=1}^N \log P(Y^{(n)} | \underline{x}, \underline{w}) = \sum_{n=1}^N [Y^{(n)} \log \sigma(\underline{w}^T \underline{x}) + (1 - Y^{(n)}) \log (1 - \sigma(\underline{w}^T \underline{x}))]$$

Noting that  $\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x)) = \sigma(x)\sigma(-x)$ ,

$$\frac{dL(\underline{w})}{d\underline{w}} = \frac{d}{d\underline{w}} \sum_{n=1}^N [Y^{(n)} \log \sigma(\underline{w}^T \underline{x}) + (1 - Y^{(n)}) \log (1 - \sigma(\underline{w}^T \underline{x}))]$$

$$= \sum_{n=1}^N \left[ \frac{\partial}{\partial \underline{w}} (Y^{(n)} \log \sigma(\underline{w}^T \underline{x})) + (1 - Y^{(n)}) \log (1 - \sigma(\underline{w}^T \underline{x})) \right]$$

$$= \sum_{n=1}^N \left[ Y^{(n)} \frac{\underline{x}}{\sigma(\underline{w}^T \underline{x})} \sigma(\underline{w}^T \underline{x}) (1 - \sigma(\underline{w}^T \underline{x})) + (1 - Y^{(n)}) \frac{\underline{x}}{\sigma(\underline{w}^T \underline{x})} \sigma(\underline{w}^T \underline{x}) (-\sigma(\underline{w}^T \underline{x})) \right]$$

$$\frac{dL(\underline{w})}{d\underline{w}} = \sum_{n=1}^N [(Y^{(n)} - \sigma(\underline{w}^T \underline{x})) \underline{x}]$$

- There is no closed form soln for the MLE  $\rightarrow$  we numerical methods (gradient ascent)

# For Personal Use Only -bkwk2

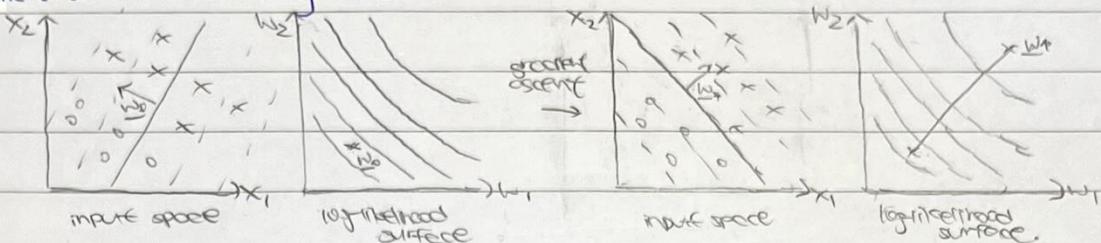
## Gradient ascent

- The batch gradient ascent rule to maximise  $L(w)$  is

$$w^{new} \leftarrow w^{old} + \alpha \frac{dL(w^{old})}{dw}$$

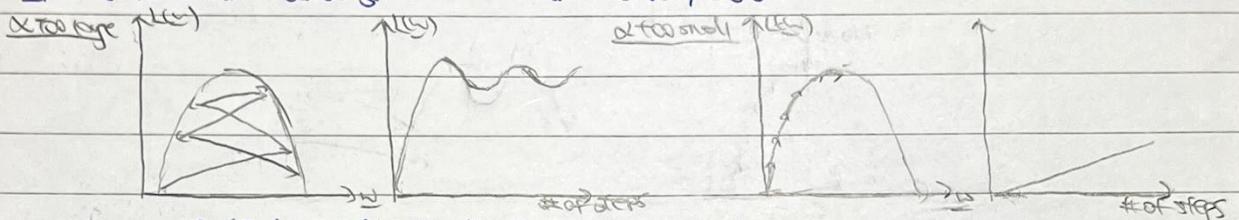
where  $\alpha > 0$  is the learning rate.

depends on  $y^{(n)} \in \{-1, 1\}$  vs  $y^{(n)} \in \{0, 1\}$ .



- If  $\alpha$  is too large  $\rightarrow$  optimisation bounces around the maximum, and could even diverge.

If  $\alpha$  is too small  $\rightarrow$  convergence to the maximum is very slow



There is no rule to choose  $\alpha$  optimally  $\rightarrow$  find  $\alpha$  by trial and error.

Linear classification with more than 2 classes.

- We can map multiple dops to discrete class labels using the max function.

$$y^{(n)} = \underset{k \in \{1, \dots, C\}}{\operatorname{argmax}} w_k^T x^{(n)}$$

but this has similar problems as in the deterministic binary classification case

- We can instead use the softmax function to map the dops into class probabilities

equivalent to  $p(y^{(n)} | x^{(n)}, w)$  for logistic regression

$$p(y^{(n)} = k | w_1, \dots, w_C, x^{(n)}) = \frac{\exp(w_k^T x^{(n)})}{\sum_{i=1}^C \exp(w_i^T x^{(n)})}$$

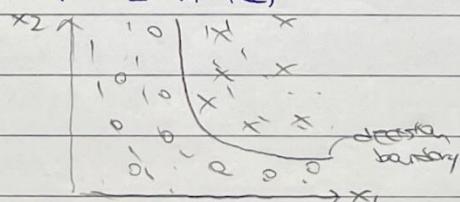
Note this is equivalent to logistic regression when  $C=2$  (softmax is a generalisation of linear regression)

Non-linear logistic regression.

- We can replace  $x$  with non-linear functions of the inputs  $\phi(x) = [\phi_1(x) \dots \phi_n(x)]^T$  to convert from linear to non-linear logistic regression.

- Non-linear logistic regression gives non-linear decision boundaries.

\* Inference does not change, we just replace  $x^{(n)}$  with  $\phi(x)$ .



## Dimensionality reduction and PCA

## Dimensionality reduction and PCA

## Dimensionality reduction

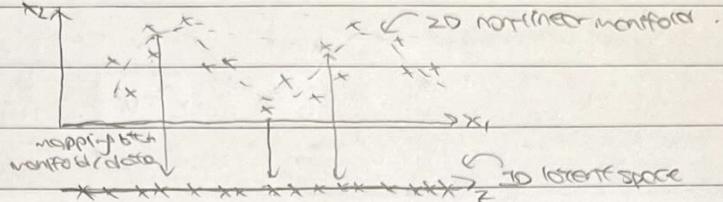
- Dimensionality reduction is a type of problem in machine learning requiring to find

↳ A low dimensional representation (manifold) of data.

↳ A mapping from data to manifold coordinates and back.

- In general, the resulting manifold is highly nonlinear, but it can be unfolded and mapped non-linearly into a latent plane

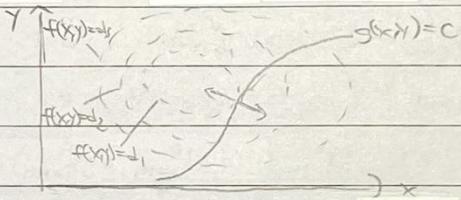
- e.g.: A 2D manifold in 2D space is unfolded into a latent real line



Interpolation b/w pts stays in the data manifold in the 1D latent space, but not in the 2D space.

## Lagrangian multipliers.

- The method of Lagrangian multipliers allows us to solve optimisation problems w/ equality constraints i.e. maximise  $f(x,y)$  subject to  $g(x,y) = c$  ↳ contour line of  $g(x,y)$



At the sol'n, the gradients of  $f(x,y)$  and  $g(x,y)$  are aligned,  $\nabla f(x,y) \parallel \nabla g(x,y)$

→ The sol'n should satisfy (i)  $g(x,y) = c$ , (ii)  $\nabla f(x,y) = -\lambda \nabla g(x,y)$

- This is equivalent to maximizing the Lagrangian function  $L(x,y,\lambda)$  wrt  $x,y,\lambda$ .

$$L(x,y,\lambda) = f(x,y) - \lambda(g(x,y) - c)$$

where  $\lambda$  is the Lagrange multiplier

\* Note that  $\nabla_{x,y} L(x,y,\lambda) = 0 \Leftrightarrow \nabla_{x,y} f(x,y) = -\lambda \nabla_{x,y} g(x,y)$ ;  $\nabla_\lambda L(x,y,\lambda) = 0 \Leftrightarrow g(x,y) = c$

- If there are multiple restrictions  $h_i(x) = 0, i=1,2\dots$ , the Lagrangian becomes  $L(x,y,\lambda)$

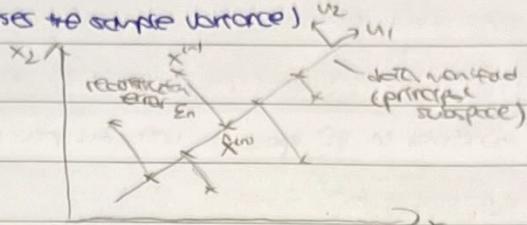
$$L(x,y,\lambda) = f(x,y) - \sum_i \lambda_i h_i(x,y)$$

# For Personal Use Only -bkwk2

## Principal Component Analysis (PCA)

- PCA is a linear dimensionality method - data manifold is assumed to be linear.
- PCA finds the projection that minimises the square reconstruction error (equivalently,

the projection that maximises the sample variance)



we can always shift  
by a const. s.t.  
this is the case

- Assume we are given a dataset  $D = \{X^{(n)}\}_{n=1}^N$  w/  $X^{(n)} \in \mathbb{R}^D$ , and  $\sum_{n=1}^N X^{(n)} = 0$

using the orthonormal basis  $\{u_i\}_{i=1}^D$ , we can rewrite each dataset  $X^{(n)}$  as

$$X^{(n)} = \underline{X}^{(n)} + \underline{\epsilon}_n = \sum_{i=1}^D \underline{X}^{(n)T} u_i u_i^T + \sum_{i=M+1}^D \underline{X}^{(n)T} u_i u_i^T$$

$u_i$ : principal comp. vectors  
 $\underline{X}^{(n)T} u_i$ : principal comp. score  
projected vector  $\underline{X}^{(n)}$  reconstruction error  $\underline{\epsilon}_n$

the latent coordinate in new basis is  $\underline{x}^{(n)} = \underline{U}^{(n)}$ , where  $\underline{U} = [\underline{u}_1 \ \underline{u}_2 \ \dots \ \underline{u}_M]$  is the  $M \times D$  projection matrix

- The average squared reconstruction errors is given by. (Note  $\underline{\epsilon}_n = \sum_{i=M+1}^D \underline{X}^{(n)T} u_i u_i^T$ )

$$\text{cost}(\underline{U}^{(D)}) = \frac{1}{N} \sum_{n=1}^N \underline{\epsilon}_n^T \underline{\epsilon}_n = \frac{1}{N} \sum_{n=1}^N \left( \sum_{i=M+1}^D \underline{X}^{(n)T} u_i u_i^T \right) \left( \sum_{j=M+1}^D \underline{X}^{(n)T} u_j u_j^T \right)$$

$$= \frac{1}{N} \sum_{n=1}^N \sum_{i=M+1}^D \sum_{j=M+1}^D \underline{u}_i^T \underline{X}^{(n)} \underline{X}^{(n)T} \underline{u}_j = \frac{1}{N} \sum_{n=1}^N \sum_{i=M+1}^D \underline{u}_i^T \underline{X}^{(n)T} \underline{X}^{(n)T} \underline{u}_j$$

orthogonal property  
 $\underline{u}_i^T \underline{u}_j = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$

$$= \sum_{i=M+1}^D \underline{u}_i^T \left( \frac{1}{N} \sum_{n=1}^N \underline{X}^{(n)T} \underline{X}^{(n)} \right) \underline{u}_i = \sum_{i=M+1}^D \underline{u}_i^T \underline{\Sigma} \underline{u}_i$$

where  $\underline{\Sigma} = \frac{1}{N} \sum_{n=1}^N \underline{X}^{(n)T} \underline{X}^{(n)}$  is the covariance matrix for the data.

- We use the method of Lagrange multipliers to min. the cost to ensure that  $\{\underline{u}_i\}_{i=1}^D$  has unit norm.

$$L(\{\underline{u}_i\}_{i=1}^D, \lambda) = \sum_{i=1}^D \underbrace{[\underline{u}_i^T \underline{\Sigma} \underline{u}_i + \lambda(1 - \underline{u}_i^T \underline{u}_i)]}_{\text{cost function}}$$

$$\nabla_{\underline{u}_i} L(\{\underline{u}_i\}_{i=1}^D, \lambda) = 2\underline{u}_i^T \underline{\Sigma} \underline{u}_i - 2\lambda \underline{u}_i = 0 \rightarrow \underline{u}_i^T \underline{\Sigma} \underline{u}_i = \lambda \underline{u}_i^T \underline{u}_i$$

Note  $\underline{u}_i^T (\underline{\Sigma} \underline{u}_i) = \underline{u}_i^T (\underline{\Sigma} \underline{u}_i^T) = 2\lambda$

$$\nabla_{\lambda} L(\{\underline{u}_i\}_{i=1}^D, \lambda) = 1 - \underline{u}_i^T \underline{u}_i = 0 \rightarrow \underline{u}_i^T \underline{u}_i = 1$$

→  $(\lambda_i, \underline{u}_i)$  are the eigenvalue and eigenvectors of the covariance matrix  $\underline{\Sigma}$ .

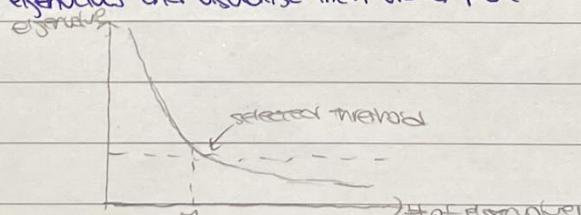
- We can write the objective function in terms of this sol'n

$$\text{cost}(\{\underline{u}_i\}_{i=1}^D) = \sum_{i=1}^D \underline{u}_i^T \underline{\Sigma} \underline{u}_i = \sum_{i=1}^D \underline{u}_i^T (\lambda_i \underline{u}_i) = \sum_{i=1}^D \lambda_i (\underline{u}_i^T \underline{u}_i) = \sum_{i=1}^D \lambda_i$$

The PCA sol'n is given by the  $\{\underline{u}_i\}_{i=1}^D$  s.t.  $\underline{u}_1, \dots, \underline{u}_D$  are eigenvectors of  $\underline{\Sigma}$  w/ the D-M

smallest eigenvalues;  $\underline{u}_1, \dots, \underline{u}_M$  are eigenvectors of  $\underline{\Sigma}$  w/ the M largest eigenvalues.

- To choose  $M$ , we sort the eigenvalues and visualise them via a plot



We choose  $M$  to be the value beyond which there is an inflection pt. The remaining eigenvalues are small w/ eigenvectors that only capture noise.

## K-means clustering

## clustering

- In general, two pts belonging to the same cluster are more similar/closer to each other than two pts belonging to different clusters.
- clustering is an unsupervised learning problem; no labels or rewards.
- clustering maps the dataset  $D = \{x^{(n)}\}_{n=1}^N$  to a vector of cluster assignments  $\mathbf{z} = \{z_1, \dots, z_N\}$

## K-means clustering

we cannot use MLE to find  $k$ ,  
or we find  $k$  to be equal to  $N$

- In K-means clustering, we need to specify the no. of clusters  $K$

- For input  $D = \{x^{(n)}\}_{n=1}^N, x^{(n)} \in \mathbb{R}^p$ , the K-means algorithm is as follows:

4 From initial (random) cluster means  $\underline{m}_k \in \mathbb{R}^p$  for  $k=1 \dots K$ , iterate the two steps

(i) Assignment step:  $s_{nk} = \arg \min_{k \in \{1, \dots, K\}} \|x^{(n)} - \underline{m}_k\|_2$  for  $n=1 \dots N$  [assign to closest cluster]

(ii) Update step:  $\underline{m}_k = \text{mean}(x^{(n)} : s_{nk}=k)$  for  $k=1 \dots K$  [update the cluster means]

- K-means clustering can be interpreted as optimisation. Let  $s_{nk} = 1$  if data pt.  $n$  is assigned to cluster  $k$ , and zero otherwise. (Note  $\sum_{k=1}^K s_{nk} = 1$ , i.e. sum of rows of  $\mathbf{s}$  equals 1, where  $s_{nk} = (\mathbf{s})_{nk}$ )

$$C(\{s_{nk}\}, \{\underline{m}_k\}) = \sum_{n=1}^N \sum_{k=1}^K s_{nk} \|x^{(n)} - \underline{m}_k\|_2^2 \quad \begin{matrix} \text{could think as PE of a system} \\ \text{of springs connecting } x^{(n)} \text{ to } \underline{m}_k \end{matrix}$$

K-means tries to minimise the cost function  $C$  wrt  $\{s_{nk}\}$  and  $\{\underline{m}_k\}$ , subject to  $\sum_{k=1}^K s_{nk} = 1$  and  $s_{nk} \in \{0, 1\}$

(i) Assignment step: minimise  $C$  wrt  $\{s_{nk}\}$ , holding  $\{\underline{m}_k\}$  fixed.

As each data pt. is independent, minimising  $C \Leftrightarrow$  minimising  $\sum_{k=1}^K s_{nk} \|x^{(n)} - \underline{m}_k\|_2^2$  for each row  $n$ .

We minimise  $\sum_{k=1}^K s_{nk} \|x^{(n)} - \underline{m}_k\|_2^2$  when  $s_{nk}=1$  for min.  $\|x^{(n)} - \underline{m}_k\|_2^2$  and zero otherwise  $\rightarrow s_n = \arg \min_k \|x^{(n)} - \underline{m}_k\|_2^2$

(ii) Update step: minimise  $C$  wrt  $\{\underline{m}_k\}$ , holding  $\{s_{nk}\}$  fixed.

As each cluster is independent, minimising  $C$  by minimising  $\sum_{n=1}^N \|x^{(n)} - \underline{m}_k\|_2^2$  for each cluster  $k$

$$\frac{\partial C}{\partial \underline{m}_k} = \sum_{n=1}^N \left[ -2x^{(n)} + 2\underline{m}_k \right] = 0 \rightarrow \underline{m}_k = \frac{\sum_{n=1}^N x^{(n)} : s_{nk}}{\sum_{n=1}^N s_{nk}} = \text{mean}(x^{(n)} : s_{nk}).$$

- The cost function  $C$  is a Lyapunov function, so K-means converges. (except for edge cases)

(finding the global optimum of the cost function  $C$  is a NP-hard problem)

- There are several issues w/ K-means clustering:

$\hookrightarrow$  The result depends on the initial cluster means  $\underline{m}_k \rightarrow$  use k-means++ algorithm to spread out the initial centres, or repeat the algorithm multiple times and pick the best one.

$\hookrightarrow$  The algorithm returns hard assignments  $\rightarrow$  use EM algorithm for soft assignments (probabilities)

$\hookrightarrow$  The algorithm assumes spherical clusters  $\rightarrow$  gives strange results when the clusters are not spherical.

# For Personal Use Only -bkwk2

## EM algorithm

### Jensen's inequality

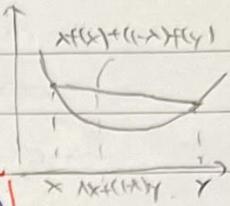
- Jensen's inequality is a general result of convexity. For a convex function  $f$  and  $\lambda \in [0, 1]$

$$\lambda f(x) + (1-\lambda)f(y) \geq f(\lambda x + (1-\lambda)y)$$

we can generalise the result to expectation

$$E[f(x)] \geq f(E[x])$$

\* The inequality is reversed for concave functions, so  $E[f(x)] \leq f(E[x])$



### Kullback-Leibler (KL) divergence

- the KL divergence b/w two probability distributions  $p(z)$  and  $p_2(z)$  is defined as

$$KL(p_1(z) || p_2(z)) = \sum_z p_1(z) \log \frac{p_1(z)}{p_2(z)}$$

If it is a measure of the "distance" / how similar two probability distributions are

- the KL divergence has the following properties:

(i) Nonsymmetric:  $KL(p_1(z) || p_2(z)) \neq KL(p_2(z) || p_1(z))$  → "divergence" not "distance"

(ii) Gibbs' inequality:  $KL(p_1(z) || p_2(z)) \geq 0$ , w/ equality iff  $p_1(z) = p_2(z)$ . ↪ proof in JFT notes

### Negative variational free energy $F(\theta, q(\cdot|z))$

- The negative variational free energy / evidence lower bound (ELBO) is a lower bound

for the likelihood  $\log p(\underline{x}|\theta)$ , and is defined to be

$$F(\theta, q(\cdot|z)) = \log p(\underline{x}|\theta) - KL(q(\cdot|z) || p(\cdot|z|\underline{x}, \theta)) = \log p(\underline{x}|\theta) - \sum_z q(z) \log \frac{p(z|\theta)}{p(z|\underline{x}, \theta)}$$

- Note that since KL-divergence is non-negative, we have  $F(\theta, q(\cdot|z)) \leq \log p(\underline{x}|\theta)$ , w/ equality

iff  $q(z) = p(z|\underline{x}, \theta)$ .

- To derive the negative variational free energy  $F(\theta, q(\cdot|z))$ , consider the log-likelihood  $\log p(\underline{x}|\theta)$

$$\begin{aligned} \log p(\underline{x}|\theta) &= \log \left[ \sum_s p(\underline{x}, s|\theta) \right] = \log \left[ \sum_s p(s|\theta) p(\underline{x}|s, \theta) \right] = \log \left[ \sum_s \frac{p(s|\theta) p(\underline{x}|s, \theta)}{q(s)} \cdot \frac{q(s)}{q(s)} \right] \\ &= \log \left[ E_q \left[ \frac{p(s|\theta) p(\underline{x}|s, \theta)}{q(s)} \right] \right] \stackrel{\text{Jensen}}{\geq} E_q \left[ \log \frac{p(s|\theta) p(\underline{x}|s, \theta)}{q(s)} \right] = F(\theta, q(\cdot|z)) \end{aligned}$$

$$\therefore F(\theta, q(\cdot|z)) = E_q \left[ \log p(s|\theta) p(\underline{x}|s, \theta) \right] - E_q \left[ \log q(s) \right] = \sum_z q(z) \log p(z|\theta) p(\underline{x}|z, \theta) - \sum_z q(z) \log q(z)$$

this form is equivalent to the form above. Noting that  $p(\underline{x}, s|\theta) = p(\underline{x}|\theta) p(s|\underline{x}, \theta) = p(s|\theta) p(\underline{x}|s, \theta)$ ,

$$\begin{aligned} F(\theta, q(\cdot|z)) &= E_q \left[ \log p(\underline{x}|\theta) p(s|\underline{x}, \theta) \right] - E_q \left[ \log q(s) \right] = \sum_z q(z) \log p(\underline{x}|\theta) p(\underline{x}|z, \theta) - \sum_z q(z) \log q(z) \\ &= \sum_z q(z) \log p(\underline{x}|\theta) + \sum_z q(z) \log \frac{p(z|\theta)}{q(z)} \\ &= \log p(\underline{x}|\theta) \sum_z q(z) - \sum_z q(z) \log \frac{p(z|\theta)}{q(z)} = \log p(\underline{x}|\theta) - KL(q(\cdot|z) || p(\cdot|z|\underline{x}, \theta)) \end{aligned}$$

+  $F(\theta, q(\cdot|z)) = \log p(\underline{x}|\theta) - KL(q(\cdot|z) || p(\cdot|z|\underline{x}, \theta))$  is useful for the E-step (determin cont. w/  $q(\cdot|z)$ )

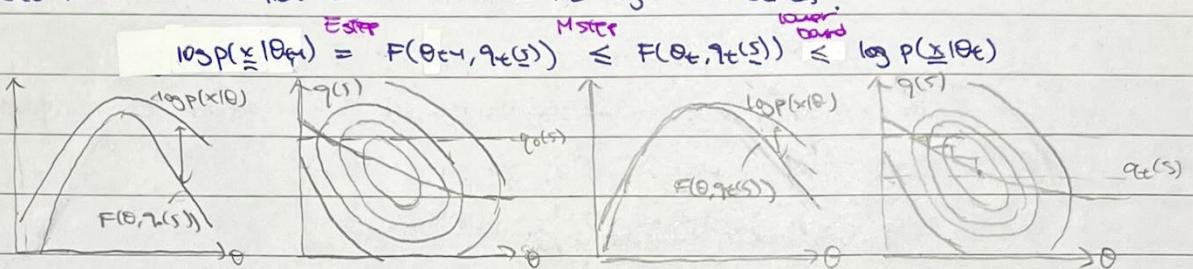
$F(\theta, q(\cdot|z)) = \sum_z q(z) \log p(z|\theta) p(\underline{x}|z, \theta) - \sum_z q(z) \log q(z)$  is useful for the M-step (one term const w/  $\theta$ )

# For Personal Use Only -bkwk2

## The expectation maximization (EM) algorithm

- The EM algorithm is an iterative method to find (local) MLE of parameters  $\theta$  in statistical models which involves unobserved latent variables  $s$ .
- Since the negative variational free energy  $F(\theta, q(s))$  is a lower bound of the log likelihood  $p(x|\theta)$ , optimising  $F(\theta, q(s))$  gives a better and better lower bound for  $\log p(x|\theta)$ .
- The EM algorithm is as follows:
  - ↳ From initial (random) parameters  $\theta_0$ , iterate  $t=1\dots T$  the two steps
  - (i) E-step: For fixed  $\theta_{t-1}$ , maximise lower bound  $F(\theta_{t-1}, q_t(s))$  wrt  $q_t(s)$
  - $$F(\theta_{t-1}, q_t(s)) = \log p(x|\theta_{t-1}) - KL(q_t(s) || p(s|x, \theta_{t-1})) \rightarrow q_t(s) = p(s|x, \theta_{t-1})$$
  - (ii) M-step: For fixed  $q_t(s)$ , maximise the lower bound  $F(\theta, q_t(s))$  wrt  $\theta$ .
  - $$F(\theta, q_t(s)) = \sum_s q_t(s) \log p(x|s, \theta) p(s|\theta) - \sum_s q_t(s) (\log q_t(s)) \rightarrow \theta_t = \arg\max_{\theta} \sum_s q_t(s) \log p(x|s, \theta) p(s|\theta)$$

- Note that each iteration cannot decrease the log likelihood as.



## Mixture of Gaussians (MoG) generative model

- In the MoG model, for each data pt  $n=1\dots N$ ,
  - ↳ the sample cluster membership distribution is given by
  - $$p(s_n=k|\theta) = \pi_k, \text{ where } \sum_{k=1}^K \pi_k = 1$$
  - ↳ the sample data value distribution given cluster membership is given by
  - $$p(x^{(n)}|s_n=k, \theta) = N(x^{(n)}|\mu_k, \Sigma_k)$$
- We would like to first learn the parameters  $\theta$  of the MoG model from data using max. likelihood,  $\theta^{ML} = \arg\max_{\theta} \log p(x|\theta)$ , then infer the cluster membership from the observed data  $p(s_n=k|x^{(n)}, \theta^{ML})$
- The log-likelihood  $\log p(x|\theta)$  is given by
 
$$\begin{aligned} \log p(x|\theta) &= \log \prod_{n=1}^N p(x^{(n)}|\theta) = \sum_{n=1}^N \log p(x^{(n)}|\theta) = \sum_{n=1}^N \log \sum_{k=1}^K p(x^{(n)}, s_n=k|\theta) \\ &= \sum_{n=1}^N \log \sum_{k=1}^K p(s_n=k|\theta) p(x^{(n)}|s_n=k, \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k N(x^{(n)}|\mu_k, \Sigma_k) \end{aligned}$$
- This can be maximised using the EM algorithm or using direct gradient-based optimisation.
- Direct optimisation is generally faster, but the EM algorithm is often simpler to implement, more widely used and can lead to important generalisations.

# For Personal Use Only -bkwk2

Application of the EM algorithm to the MoG model.

- Assume the data is 1D ( $x \in \mathbb{R}$ ) for simplicity. Let the Gaussian mixture model parameters be  $\theta = \{\mu_k, \sigma_k^2, \pi_k\}_{k=1}^K$  and 1 latent variable per data pt  $S_n \in \{1 \dots K\}$  for  $n=1 \dots N$

- Applying the EM algorithm to the simplified MoG model,

(i) E-step: Note that  $p(x^{(n)} | S_n=k, \theta) = \frac{1}{2\pi\sigma_k^2} \exp\left[-\frac{1}{2\sigma_k^2}(x^{(n)} - \mu_k)^2\right]$ ,  $P(S_n=k|\theta) = \pi_k$

$$q(S_n=k) = P(S_n=k | X^{(n)}, \theta) \propto p(X^{(n)} | S_n=k, \theta) p(S_n=k | \theta) = \frac{\pi_k}{2\pi\sigma_k^2} \exp\left[-\frac{1}{2\sigma_k^2}(x^{(n)} - \mu_k)^2\right] = \text{unk}$$

$$\therefore q(S_n=k) = \frac{\pi_k}{\sum_{j=1}^K \pi_j} = r_{n,k} \quad (\text{"responsibility" that component } k \text{ takes for data pt } n)$$

(ii) M-step:  $F(\theta, q(S)) = \sum_S q(S) \log p(S | \theta) p(X | S, \theta) - \sum_S q(S) \log q(S)$

$$= \sum_{n=1}^N \sum_{k=1}^K q(S_n=k) [\log \pi_k - \frac{1}{2} \log \sigma_k^2 - \frac{1}{2\sigma_k^2} (x^{(n)} - \mu_k)^2] + \text{const.}$$

$$\therefore \frac{\partial F}{\partial \mu_j} = \sum_{n=1}^N q(S_n=j) \frac{x^{(n)} - \mu_j}{\sigma_j^2} = 0 \rightarrow \mu_j = \frac{\sum_{n=1}^N q(S_n=j) x^{(n)}}{\sum_{n=1}^N q(S_n=j)}$$

$$\frac{\partial F}{\partial \sigma_j^2} = \sum_{n=1}^N q(S_n=j) \left[ -\frac{1}{2\sigma_j^2} + \frac{(x^{(n)} - \mu_j)^2}{2\sigma_j^4} \right] = 0 \rightarrow \sigma_j^2 = \frac{\sum_{n=1}^N q(S_n=j) (x^{(n)} - \mu_j)^2}{\sum_{n=1}^N q(S_n=j)}$$

$$\frac{\partial (F + \lambda(1 - \sum_k \pi_k))}{\partial \pi_j} = \sum_{n=1}^N \frac{q(S_n=j)}{\pi_j} - \lambda = 0 \rightarrow \pi_j = \frac{1}{N} \sum_{n=1}^N q(S_n=j) \quad \begin{cases} \sum_{j=1}^K q(S_n=j) = 1 \\ \sum_{n=1}^N q(S_n=j) = N \end{cases}$$

- In summary, the EM algorithm for the MoG model is as follows:

↳ From initial (random) parameters  $\theta = \{\mu_k, \sigma_k^2, \pi_k\}_{k=1}^K$ , iterate the two steps:

(i) E-step:  $r_{n,k} = p(S_n=k | X^{(n)}, \theta)$  for  $n=1 \dots N$

(ii) M-step:  $\theta = \arg \max \sum_{n=1}^N \sum_{k=1}^K r_{n,k} p(X^{(n)} | \theta, S_n=k)$

- Applying the EM algorithm to the MoG model is soft k-means clustering w/ non-axix aligned, non-equally weighted clusters.

- Applying the EM algorithm to the MoG model has several drawbacks:

↳ MoG clusters still have simple shapes → one cluster can have many components / use complex models

↳ ML approach could overfit → use Bayesian approach

↳ coordinate descent is often slow to converge → joint optimisation of  $F(\theta, q(S))$  or direct optimisation of the log-likelihood  $\log(p(X | \theta))$

↳ The free energy may diverge (as it is not bounded above) → "koom problem" when one of the data pts take the entire cluster.

# For Personal Use Only -bkwk2

sequence modelling

## MARKOV models

sequence modelling.

- sequence modelling can be used to:

↳ predict future items in a sequence  $p(Y_{t+1}|Y_1, \dots, Y_t)$  [e.g.: FTSE 100 trends]

↳ remove noise from a sequence  $p(Y_{t+1}|Y_1, \dots, Y_t)$  [e.g.: Noise removal]

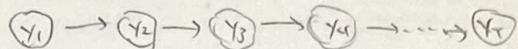
↳ predict one sequence from another  $p(Y_{t+1}|Y_1, \dots, Y_t)$  [e.g.: translation]

↳ discover underlying latent variables  $p(X_{t+1}|Y_1, \dots, Y_t)$  [e.g.: sentiment analysis]

Markov models.

- A Markov model has the Markov (memoryless) property.

(i) 1st order Markov:  $Y_{t+1} \perp Y_1, \dots, Y_t | Y_t$   $[p(Y_{t+1}|Y_t) = p(Y_1)p(Y_2|Y_1)p(Y_3|Y_2)\dots p(Y_t|Y_{t-1})]$



(ii) 2nd order Markov:  $Y_{t+1} \perp Y_1, \dots, Y_{t-2} | Y_{t-1}, Y_t$   $[p(Y_{t+1}|Y_{t-1}, Y_t) = p(Y_1)p(Y_2|Y_1)p(Y_3|Y_1, Y_2)\dots p(Y_t|Y_{t-2}, Y_{t-1})]$



- Markov models for discrete data — N-gram models

(i) 1st order Markov (Bigram):  $Y_t \in \{1, \dots, K\}$ ,  $p(Y_1=k) = \pi_k^0$ ,  $p(Y_t=k|Y_{t-1}=l) = T_{k,l}$

(ii) 2nd order Markov (Trigram):  $Y_t \in \{1, \dots, K\}$ ,  $p(Y_1=k) = \pi_k^0$ ,  $p(Y_t=k|Y_{t-1}=l, Y_{t-2}=m) = T_{k,l,m}$

- For the bigram model, the marginal distribution for the second state  $p(Y_2=k)$  is given by

$$p(Y_2=k) = \sum_{l=1}^K p(Y_2=k|Y_1=l) p(Y_1=l) = \sum_{l=1}^K T_{k,l} \pi_l^0 = [\Pi \Pi^0]_k \rightarrow p(Y_2) = \Pi \Pi^0$$

We can find the stationary distribution of the Markov chain by finding the eigenvectors of

the transition matrix  $\Pi$  w/ eigenvalue equal to one, since

$$p(Y_t=k) = \sum_{l=1}^K p(Y_t=k|Y_{t-1}=l) p(Y_{t-1}=l) \rightarrow \Pi^{\infty} = \sum_{l=1}^K T_{k,l} \Pi^0 \rightarrow \Pi^{\infty} = \Pi \Pi^{\infty} \quad (\lambda=1)$$

\* We can find the max. likelihood transition matrix from an observed sequence by considering the proportions of transitions from state  $i$  to state  $j$  (out of all transitions from state  $i$ )

- Markov models for continuous data — Auto-regressive (AR) models

(i) 1st order Markov (AR(1)):  $Y_t \in \mathbb{R}^p$ ,  $p(Y_t) = N(Y_t | M_0, \Sigma_0)$ ,  $p(Y_t|Y_{t-1}) = N(Y_t | \Delta Y_{t-1}, \Sigma)$  (MAG)

(ii) 2nd order Markov (AR(2)):  $Y_t \in \mathbb{R}^p$ ,  $p(Y_t) = N(Y_t | M_0, \Sigma_0)$ ,  $p(Y_t|Y_{t-1}, Y_{t-2}) = N(Y_t | \Delta_1 Y_{t-1} + \Delta_2 Y_{t-2}, \Sigma)$

- For AR(1) and  $Y_t \in \mathbb{R}^p$ ,  $p(Y_t|Y_{t-1}) = N(Y_t | \lambda Y_{t-1}, \sigma^2) \rightarrow Y_t = \lambda Y_{t-1} + \sigma \varepsilon_t$ ,  $\varepsilon_t \sim N(0, 1)$

Since  $p(Y_{t-1})$  is Gaussian, and the marginal of a Gaussian is Gaussian, the stationary distribution  $p(Y_0)$  is Gaussian

↳ mean:  $E[Y_t] = \lambda E[Y_{t-1}] + \sigma E[\varepsilon_t] = \lambda E[Y_{t-1}] \rightarrow M_0 = \lambda M_{t-1}$

$$\text{As } t \rightarrow \infty, M_{t-1} = \lambda M_0 \rightarrow M_0 = 0$$

↳ variance:  $E[Y_t^2] = \lambda^2 E[Y_{t-1}^2] + 2\lambda \sigma E[Y_{t-1} \varepsilon_t] + \sigma^2 E[\varepsilon_t^2] = \lambda^2 E[Y_{t-1}^2] + 2\lambda \sigma E[Y_{t-1}] E[\varepsilon_t] + \sigma^2 \rightarrow \sigma_0^2 = \lambda^2 \sigma_0^2 + \sigma^2$

$$\text{As } t \rightarrow \infty, \sigma_0^2 = \lambda \sigma_0^2 + \sigma^2 \rightarrow \sigma_0^2 = \frac{\sigma^2}{1-\lambda^2}$$

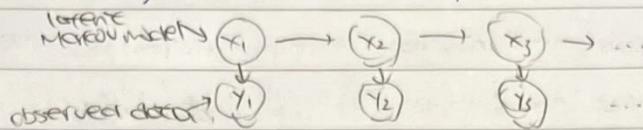
+ We always set  $\Delta$  s.t. its eigenvalue has absolute value less than 1  $\rightarrow \sigma_0^2$  does not blow up

# For Personal Use Only -bkwk2

## Hidden Markov models

### Hidden Markov models

- In hidden Markov models, the real (observed) data  $y_t$  depend on latent (unobserved) variables  $x_t$ .



The joint distribution over all the latent variables  $x_{1:T}$  and real data  $y_{1:T}$  is given by

$$p(y_{1:T}, x_{1:T}) = \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t), \text{ where } p(x_t | x_{t-1}) = p(x_t)$$

We can marginalise the latents and get a fully-connected graph (not HMM now)



- Two prevalent examples include:

↳ Hidden Markov model (HMM): discrete  $x$ , discrete/continuous  $y$

↳ Linear Gaussian state space model (LGSSM): continuous  $x$ , continuous  $y$ .

### Hidden Markov model — discrete hidden state (HMM)

- HMMs have discrete hidden states:  $x_t \in \{1 \dots K\}$ ,  $p(x_t=k | x_{t-1}=l) = T_{k,l}$ .

- The observed state can be continuous or discrete

↳ Continuous observed state:  $p(y_t | x_t=k) = N(y_t | \mu_k, \Sigma_k)$ .

↳ Discrete observed state:  $p(y_t=l | x_t=k) = \pi_{k,l}$ , where  $\pi$  is the emission matrix.

- For discrete hidden state, continuous observed state,

$$p(y_t) = \sum_{k=1}^K p(y_t | x_t=k) p(x_t=k) = \sum_{k=1}^K N(y_t | \mu_k, \Sigma_k) T_{k,k}$$

$$p(y_{1:T}) = \sum_{k=1}^K p(y_{1:T} | x_{1:T}=k) p(x_{1:T}=k) = \sum_{k=1}^K N(y_{1:T} | \mu_k, \Sigma_k) T_{k,k}^{\otimes T}$$

→ we can interpret them as a MoG model w/ dynamic cluster assignments.

### Hidden Markov model — continuous hidden state (LGSSM)

- LGSSMs have continuous hidden state:  $x_t \in \mathbb{R}^d$ ,  $p(x_t | x_{t-1}) = N(x_t | \mu_{t-1}, Q_t)$

- The observed state is continuous

↳ Continuous observed state:  $y_t \in \mathbb{R}^d$ ,  $p(y_t | x_t) = N(y_t | \mu_{x_t}, R_t)$

# For Personal Use Only -bkwk2

## Varieties of Inference

- The four distributional estimates are as follows:

	Interstate score vs sequence? (single state) marginal	Sequence joint
filter online	$p(x_t   y_{1:t})$	$p(x_{1:t}   y_{1:t})$
smoother offline	$p(x_t   y_{1:T})$	$p(x_{1:t}   y_{1:T})$

A diff. algorithm is req. for each of the distributional estimates  $\rightarrow$  JFG focuses on  $p(x_t | y_{1:t})$

- The two point estimates are as follows:

$$\hookrightarrow x_t^* = \underset{x_t}{\operatorname{argmax}} p(x_t | y_{1:t}) \quad \text{"bit-optimal"} \quad \hookrightarrow x_{1:T}^* = \underset{x_{1:T}}{\operatorname{argmax}} p(x_{1:T} | y_{1:T}) \quad \text{"block-optimal"}$$

- The estimates  $x_{1:T}^*$  and  $x_t^*$  are the same for the LGSSM (but not discrete hidden state HMM)!

For LGSSM,  $p(x_{1:T} | y_{1:T}) = N(x_{1:T} | M_{1:T}, \Sigma_{1:T|1:T})$  [Gaussian since  $p(x_{1:T}, y_{1:T})$  is Gaussian]

$$\hookrightarrow x_t^* = \underset{x_t}{\operatorname{argmax}} p(x_t | y_{1:T}) = M_t \rightarrow x_{1:T}^* = M_{1:T}$$

$$\hookrightarrow \Sigma_{1:T}^* = \underset{\Sigma_{1:T}}{\operatorname{argmax}} p(\Sigma_{1:T} | y_{1:T}) = M_{1:T}$$

shown w/ counterexample

## Kalman Filter

- The Kalman filter is an iterative algorithm to compute  $p(x_t | y_{1:t})$  for the LGSSM.

$$(i) \text{ Start w/ } p(x_{1:t} | y_{1:t-1}) = N(x_{1:t} | M_{1:t}, V_{1:t}) \quad \text{use } x_t = Ax_{t-1} + Bv_t, v_t \sim N(0, I)$$

$$(ii) \text{ Diffuse via dynamics: } p(x_t | y_{1:t}) = \int p(x_t | x_{t-1}) p(x_{t-1} | y_{1:t-1}) dx_{t-1}$$

$$p(x_t | y_{1:t}) = \int N(x_t | Ax_{t-1}, Q) N(x_{t-1} | M_{1:t}, V_{1:t}) dx_{t-1} = N(x_t | M_{1:t}, V_{1:t})$$

$$\text{where } M_{1:t}^{*t} = A M_{1:t-1}^{*t-1}, \quad V_{1:t}^{*t} = A V_{1:t-1}^{*t-1} A^T + Q$$

$$(iii) Combine w/ likelihood: \quad p(x_t | y_{1:t}) \propto p(x_t | y_{1:t-1}) p(y_t | x_t)$$

$$\therefore p(x_t | y_{1:t}) = N(x_t | M_t, V_t)$$

$$\text{where } M_t = M_{t-1} + K_t (y_t - C M_{t-1}), \quad V_t = V_{t-1} - K_t C V_{t-1} C^T, \quad K_t = V_{t-1}^{-1} C^T (C V_{t-1}^{-1} C^T + R)^{-1}$$

## Forward algorithm

Markov property allows us to use dynamic programming / belief propagation for efficient computation.

- The forward algorithm is an iterative algorithm to compute  $p(x_t | y_{1:t})$  for discrete hidden state HMM.

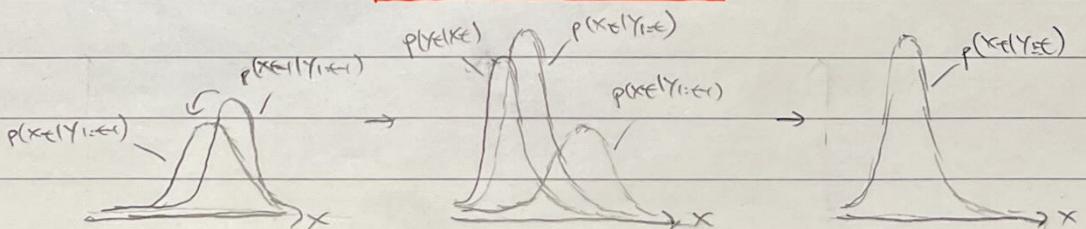
$$(i) \text{ Start w/ } p(x_{1:t} = k_t | y_{1:t-1}) = p_{t-1}^{k_{t-1}}(k_t)$$

$$(ii) \text{ Diffuse via dynamics: } p(x_{t+1} = k_{t+1} | y_{1:t}) = \sum_{l=1}^K p(x_{t+1} = k_{t+1} | x_t = l) p(x_t = l | y_{1:t})$$

$$p_t^{k_t}(k_{t+1}) = \sum_{l=1}^K T_{kl} p_{t-1}^{k_{t-1}}(l)$$

$$(iii) Combine w/ likelihood: \quad p(x_{t+1} | y_{1:t}) \propto p(x_{t+1} | y_{1:t}) p(y_t | x_t = k_t)$$

$$p_t^{k_t}(k_{t+1}) \propto p_t^{k_t}(k_{t+1}) p(y_t | x_t = k_t)$$



# For Personal Use Only -bkwk2

Derivation of the general filtering equations.

$$\textcircled{1} \quad p(x_t | Y_{1:t-1}) = \int p(x_t | x_{t-1}) p(x_{t-1} | Y_{1:t-1}) dx_{t-1} \quad \checkmark \text{ replace integrals w/ sums for the discrete case}$$

$$\hookrightarrow p(x_t | Y_{1:t-1}) = \int p(x_t, x_{t-1} | Y_{1:t-1}) dx_{t-1} = \int p(x_t | x_{t-1}, Y_{1:t-1}) p(x_{t-1} | Y_{1:t-1}) dx_{t-1}$$

$$= \int p(x_t | x_{t-1}) p(x_{t-1} | Y_{1:t-1}) dx_{t-1} \quad \checkmark \text{cond. indep. from model}$$

$$\textcircled{2} \quad p(x_t | Y_{1:t}) \propto p(y_t | x_t) p(x_t | Y_{1:t-1}) \quad \text{cond. indep. from model}$$

$$\hookrightarrow p(x_t | Y_{1:t}) = \frac{p(y_t | x_t, Y_{1:t-1}) p(x_t | Y_{1:t-1})}{p(y_t | Y_{1:t-1})} = \frac{p(y_t | x_t) p(x_t | Y_{1:t-1})}{p(y_t | Y_{1:t-1})}$$

$$\propto p(y_t | x_t) p(x_t | Y_{1:t-1})$$

Computing the likelihood.

- The likelihood  $p(Y_{1:T})$  is given by

$$p(Y_{1:T}) \prod_{t=1}^T p(y_t | Y_{1:t-1})$$

where the term  $p(y_t | Y_{1:t-1})$  is the normalisation const. for computing  $p(x_t | Y_{1:t})$

→ we simply store  $p(y_t | Y_{1:t-1})$  for each  $t$  and find the product to get the likelihood  $p(Y_{1:T})$