

## Polynomial Curve Fitting

error function:  $E(\vec{w}) = \frac{1}{2} \sum_{n=1}^N \|y(x_n - \vec{w}) - t_n\|^2$   
modified  $E(\vec{w}) = \frac{1}{2} \sum_{n=1}^N \|y(x_n - \vec{w}) - t_n\|^2 + \frac{\lambda}{2} \|\vec{w}\|^2$

$$\|\vec{w}\|^2 = \vec{w}^T \vec{w} = w_0^2 + \dots + w_N^2$$

The Rules of Probability

sum rule:  $p(X) = \sum p(X, Y)$  product rule:  $p(X, Y) = p(Y|X)p(X)$

$$\text{def: } \int_{-\infty}^{\infty} p(x) dx = 1, \quad \text{CDF: } \int_{-\infty}^x p(x) dx = P(x)$$

expectation:  $E[f] = \int p(x)f(x) dx$  or  $E[f] = \int p(x)f(x) dx$

conditional expectation:  $E_x[f] = \sum p(x,y)f(x)$

Variance:  $\text{Var}[f] = E[f(x) - E[x]]^2 = E(x^2) - [E(x)]^2$

Covariance:  $\text{cov}(x, y) = E_{xy}[(x - E[x])(y - E[y])] = E_{xy}[xy] - E[x]E[y]$

matrix:  $\text{cov}(x, y) = E_{xy}[x^T y^T] - E[x]E[y^T]$

Bayesian Prob:  $p(\vec{w}|D) = \frac{p(D|\vec{w})p(\vec{w})}{\text{prior}} \propto p(\vec{w})$

posterior  $\sqrt{\text{likelihood}} = \text{posterior} \propto \text{likelihood} \propto \text{prior}$

(normal/Gaussian)  $N(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}\right), \beta = \frac{1}{\sigma^2}, \text{precision}$

ml:  $\ln(p(x|\mu, \sigma^2)) = -\frac{1}{2} \sum_{n=1}^N (x_n - \mu)^2 \frac{1}{2\sigma^2} \ln(2\pi) - \frac{N}{2} \ln(\sigma^2)$

$\sigma_m^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_m)^2, E(\mu_m) = \mu, E(\sigma_m^2) = (\frac{N-1}{N})\sigma^2$

$\hat{\sigma}^2 = \frac{N}{N-1} \cdot \sigma_m^2 = \frac{1}{N-1} \sum (x_n - \mu_m)^2$  unbiased.

overall:  $f(x|\mu) = N^x(1-\mu)^{1-x}, E(X) = \mu, V(X) = \mu(1-\mu)$

$Mu = \frac{1}{N} \sum_{n=1}^N x_n$  sample mean

Binomial:  $Bin(m|N, p) = \binom{N}{m} (p^m)(1-p)^{N-m}, (N) = \frac{N!}{(N-m)!m!}$

Ex(m) =  $\sum_{m=0}^N m \text{Bin}(m|N, p) = Np, V(m) = \sum_{m=0}^N (m-E(m))^2 \text{Bin}(m|N, p) = N(1-p)p$

beta:  $Beta(\alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1}(1-p)^{\beta-1}, E(p) = \frac{\alpha}{\alpha+\beta}, V(p) = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$

multinomial: Bern

$$p(x^k|\vec{p}) = \prod_{k=1}^K p_k^{x_k}, \sum p(x|\vec{p}) = \sum_K p_k^{x_k}, E[x_k] = \vec{p}$$

log-likelihood:  $p(D|\vec{p}) = \prod_{n=1}^N \prod_{k=1}^K p_k^{x_{nk}} = \prod_{k=1}^K p_k^{x_{nk}}$

exp family:  $p(x|\vec{\eta}) = h(x)g(\vec{\eta}) \exp\{\vec{\eta}^T \vec{u}(x)\}, \vec{x}$  scale/vector,  $(*) (*)$

cts/dicrete.  $\vec{\eta}$  natural para of distribution  $\vec{u}(x)$

some function of  $\vec{x}$ .

Bern:  $p(x|\mu) = Bern(x|\mu) = \mu^x(1-\mu)^{1-x}, p(x|\mu) = \exp[x\ln(\mu) + (1-x)\ln(1-\mu)]$

so  $\eta = \ln(\frac{\mu}{1-\mu}), \partial_\eta \mu = \frac{1}{1+\exp(-\eta)}$

$p(x|\eta) = \mu(1-\mu) \exp(\eta x) \Rightarrow u(x) = x, h(x) = 1, g(\eta) = \mu(1-\mu)$

Multinomial:  $p(x|\vec{p}) = \vec{p}^T \vec{p}$

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

$\vec{\eta} = (\mu/\sigma^2, -\frac{1}{2\sigma^2})^T, \vec{u}(x) = (x, x^2)^T, h(\vec{\eta}) = (2\sigma)^{-\frac{1}{2}}, g(\vec{\eta}) = (2\sigma)^{\frac{1}{2}} \exp\left(\frac{\eta_1}{4\sigma^2}\right)$

ML & SS (exp family): gradient  $(*) (*)$

$\nabla g(\vec{\eta}) \int h(\vec{\eta}) \exp\{\vec{\eta}^T \vec{u}(x)\} d\vec{\eta} + g(\vec{\eta}) \int h(\vec{\eta}) \exp\{\vec{\eta}^T \vec{u}(x)\} \vec{u}(x) d\vec{\eta} = 0$

se (\*):  $-\frac{1}{\sigma^2} \nabla g(\vec{\eta}) = g(\vec{\eta}) \int h(\vec{\eta}) \exp\{\vec{\eta}^T \vec{u}(x)\} \vec{u}(x) d\vec{\eta} = E(\vec{u}(x))$

$\rightarrow \nabla g(\vec{\eta}) = E(\vec{u}(x))$

\* Likelihood function:  $p(X|\vec{\eta}) = \prod_{n=1}^N h(x_n) g(\vec{\eta})^N \exp\{\vec{\eta}^T \sum_{n=1}^N \vec{u}(x_n)\}$

$X = \{x_1, x_2, \dots, x_N\}$ , set log gradient of  $\ln(p(X|\vec{\eta})) = 0$ .

$-\nabla \ln(g(\vec{\eta})) = \frac{1}{N} \sum_{n=1}^N \vec{u}(x_n) \Rightarrow \vec{\eta}_{ML} = ?$  depends only on

$\sum_n \vec{u}(x_n) \rightarrow$  sufficient statistics

Decision theory

minimizing misclassification rate:  $p(\text{mistake}) = p(\vec{x} \in R_1, C_2) + p(\vec{x} \in R_2, C_1)$

$= \int_{R_1} p(\vec{x}, C_2) d\vec{x} + \int_{R_2} p(\vec{x}, C_1) d\vec{x}$

$p(x, C_1) p(\text{correct}) = \sum_{k=1}^K p(\vec{x} \in R_k, C_k) = \sum_{k=1}^K \int_{R_k} p(\vec{x}, C_k) d\vec{x}$

minimizing expected loss:  $E[L] = \sum_k \sum_j \int_{R_j} L_j p(\vec{x}, C_k) d\vec{x}$

loss func for reg:  $E[L] = \int_{\mathbb{R}} L(t, y(x)) p(\vec{x}, t) d\vec{x} dt$

$L(t, y(x)) = (y(x) - t)^2$  common,

normalized:  $\int p(x) dx = 1$

$y(\vec{x}) = f(p(t|\vec{x})) dt = E(t|\vec{x})$

bias parameter,  $\phi_j(x)$ : basis function

so, define  $\phi_j(x) = 1 \Rightarrow y(x, w) = \sum_{j=0}^M w_j \phi_j(x) = w^T \phi(x)$

maximum likelihood & least square

$\ln p(\vec{t}|\vec{w}, \rho) = \sum_{n=1}^N \ln \frac{1}{\sqrt{2\pi\rho}} \exp\left(-\frac{1}{2\rho} (t_n - \vec{w}^T \phi(x_n))^2\right)$

$E_D(\vec{w}) = \sum_{n=1}^N (t_n - \vec{w}^T \phi(x_n))^2$

$= \frac{1}{2} \sum_{n=1}^N (t_n - \vec{w}^T \phi(x_n))^2 + \frac{\lambda}{2} \|\vec{w}\|^2$

where sum of squares error function is

$E_D(\vec{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \vec{w}^T \phi(x_n))^2$

$\vec{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \vec{t}$

regularized least square: introduce to control over-fitting in the form:

$E_D(\vec{w}) + \lambda E_w(\vec{w})$

$= \frac{1}{2} \sum_{n=1}^N (t_n - \vec{w}^T \phi(x_n))^2 + \frac{\lambda}{2} \|\vec{w}\|^2$

$\lambda$ : regularization coefficient.

setting gradient of  $\vec{w}$  to be zero, get

$\vec{w} = (\lambda I + \Phi^T \Phi)^{-1} \Phi^T \vec{t}$

(ridge regression)

bias-var decomposition

expected loss = bias + variance + noise

(bias)<sup>2</sup> =  $\int E_D(y(x; D)) - h(x)^2 p(x) dx$

variance =  $\int E_D[y(x; D) - E_D(y(x; D))]^2 p(x) dx$

noise =  $\int [h(x) - t]^2 p(x) dx$

"trade-off" between bias & var.

Bayesian inference:

log of posterior = log likelihood + log of prior

$\ln p(\vec{w}|D) = -\frac{1}{2} \sum_{n=1}^N (t_n - \vec{w}^T \phi(x_n))^2 - \frac{\lambda}{2} \|\vec{w}\|^2 + \text{const.}$

predictive distribution:  $\Rightarrow$  integrating over posterior

conjugate prior:  $p(w|D) \propto N(w|m_0, S_0^{-1})$

posterior takes  $S_N(S_0^{-1} m_0 + \beta \Phi^T \Phi)^{-1}, S_N^{-1} = S_0^{-1} + \beta \Phi^T \Phi$

$w_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \vec{t} \Rightarrow w_{ML} = S_N(S_0^{-1} m_0 + \beta \Phi^T \Phi)^{-1} \vec{t}$

Bayesian model comparison

$\ln p(D) \approx \ln p(D|w_{ML}) + M \ln(\frac{dw}{d\vec{w}} \text{posterior})$

marginal likelihood function

$p(\vec{t}|\vec{w}, \beta) = \int p(\vec{t}|\vec{w}, \beta) p(\vec{w}|D) d\vec{w}$

weight parameter  $\vec{w}$

discriminant function:  $y(x) = \vec{w}^T \vec{x} + w_0$

$w$ : weight vector,  $w_0$ : bias.  $C_1$  if  $y(x) > 0$ ,  $C_2$  if  $y(x) < 0$ .  $y(x) = 0$  decision surface.

$\|\vec{w}\| = \sqrt{\vec{w}^T \vec{w}}$  normal distance from origin to surface

single solution for k-class discriminant

$y_k(x) = \vec{w}_k^T \vec{x} + w_{k,0}, k=1, \dots, K$

assign  $x$  to  $C_k$  if  $y_k(x) > y_j(x), \forall j \neq k$ .

A, B in  $R_k$ :  $y_k(A) > y_j(A)$

$y_k(B) > y_j(B)$

$\Rightarrow \alpha > 0, y_k(A) + (1-\alpha)y_B > y_j(A) + (1-\alpha)y_B$

Fisher's linear discriminant

Gaussian example (prob. generative models)

$p(x|\vec{w}_k, \Sigma_k) = \frac{1}{\sqrt{2\pi|\Sigma_k|}} \exp\left(-\frac{1}{2} (\vec{x} - \vec{w}_k)^T \Sigma_k^{-1} (\vec{x} - \vec{w}_k)\right)$

2 classes, posterior is logistic func.

$p(C_1|x) = \sigma(\vec{w}^T \vec{x} + w_0)$

where  $\vec{w} = \Sigma^{-1}(\vec{w}_1 - \vec{w}_2)$

$w_0 = -\frac{1}{2} \vec{w}_1^T \vec{w}_1 + \frac{1}{2} \vec{w}_2^T \vec{w}_2 + \ln \frac{p(C_1)}{p(C_2)}$

reg model with square loss func.

$E[L] = \int (y(x) - t)^2 p(x, t) dx dt$

Show  $y(x)$  for which  $E[L]$  by conditional expectation

$E[t|x] = \int t p(t|x) dt$ .

$E[L] = \int \int (y(x) - t)^2 p(x, t) dt dx$

$+ \int t^2 p(x, t) dt dx - 2 \int \int y(x) p(x, t) dt dx dt$

$= \int y(x)^2 p(x) dx - 2 \int y(x) E[t|x] p(x) dx$

$= E(t^2) + E[y(x)^2] - 2 E[y(x)] E[t|x]$

$p(w|D) = \exp\left(-\frac{1}{2} \sum_{n=1}^N (t_n - \vec{w}^T \phi(x_n))^2\right) \exp\left(-\frac{\lambda}{2} \|\vec{w}\|^2\right)$

linear  $\log p(t|x, w) = -\frac{1}{2} \sum_{n=1}^N (t_n - \vec{w}^T \phi(x_n))^2 - \frac{\lambda}{2} w^T w + \text{constant}$

basis func model:  $\Rightarrow \min \frac{\ell}{2} \sum_{n=1}^N \dots + \frac{\lambda}{2} \|\vec{w}\|^2$

quadratic term:  $\lambda = \frac{\alpha}{\beta} = \alpha \sigma^2$

Bayesian logistic reg for binary classification.

$\log p(w|x, t) = \log p(t|x, w) + \log p(w|D) + \text{constant}$

$= -\frac{1}{2} w^T w + \sum_{n=1}^N [t_n \ln y_n + (1-t_n) \ln (1-y_n)] + \text{const}$

$\Rightarrow p(w|x, t) \propto \exp\left(-\frac{1}{2} w^T w\right) \cdot \prod_{n=1}^N y_n^{t_n} (1-y_n)^{1-t_n}$

where  $y_n = \sigma(w^T x_n) = \frac{1}{1 + \exp(-w^T x_n)}$

$\log p(w|x, t) \Leftrightarrow -\sum_{n=1}^N [t_n \ln y_n + (1-t_n) \ln (1-y_n)] + \frac{\lambda}{2} w^T w$

$\lambda = d, E_n = \begin{cases} -\ln y_n & \text{if } t_n = 1 \\ -\ln (1-y_n) & \text{if } t_n = 0 \end{cases}$

$\sum E_n$  is cross-entropy error in total.

derivative of  $E_n$  (crossentropy)

$E[L] = -\sum_{n=1}^N [t_n \ln y_n + (1-t_n) \ln (1-y_n)] + \frac{\lambda}{2} w^T w$

$= \sum_{n=1}^N E_n + \frac{\lambda}{2} w^T w$

$dE_n = \frac{Y_n - t_n}{Y_n(1-Y_n)}, \frac{dy_n}{dw} = Y_n(1-Y_n)x_n$

Since  $\frac{dy_n}{dw} = \frac{dx_n}{dw} = \frac{w^T \phi(x_n)}{w^T \phi(x_n)}$

$\Rightarrow \frac{dE_n}{dw} = (Y_n - t_n)x_n \Rightarrow \nabla w E[L] = \sum_{n=1}^N x_n$

Bayesian linear regression:

likelihood:  $p(t|x, w, \beta) = N(t_n | w^T \phi(x_n), \beta^2)$

conjugate prior:  $p(w|D) \propto N(w|m_0, S_0^{-1})$

posterior takes  $S_N(S_0^{-1} m_0 + \beta \Phi^T \Phi)^{-1}, S_N^{-1} = S_0^{-1} + \beta \Phi^T \Phi$

$w_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \vec{t}$

Bayesian linear regression: integrating over posterior

$p(t|x, \vec{x}, \vec{w}, \beta) = \int p(t|x, \vec{w}, \beta) p(w|D) dw$

$= N(t | m_N \phi(x), \sigma_N^2(x))$

where  $\sigma_N^2(x) = \frac{1}{S} + \phi(x)^T S_N \phi(x), m_N = \beta S_N \Phi^T \vec{x}$

$S_N = dI + \beta \Phi^T \Phi$

evidence approximation: fully Bayesian predictive distribution:  $p(t^*|x^*, w, \beta) = \int p(t^*|x^*, w, \beta) p(w|D, \alpha, \beta) dw d\alpha d\beta$

LS for classification:  $y_k(x) = x^T w_k w_k^T x, k=1, \dots, K$

$\text{or } Y(x) = \vec{w}^T \vec{x}$

$\vec{w}_k = (w_{k,0}, w_k^T)^T$

posterior:

$P(C_i|x) = \frac{P(x|C_i)P(C_i)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)} = \frac{1}{1 + \exp(-\alpha)}$

where  $\alpha = \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)} = \ln \frac{P(C_1|x)}{1 - P(C_1|x)}$

$\sigma(C_i) = 1 - \alpha(C_i)$

$\frac{d}{da} \sigma(a) = \sigma(a)(1-\sigma(a))$

Sigmoid

$\frac{\partial \ell}{\partial w} = -\sum_{n=1}^N r_n (t_n - \vec{w}^T \phi(x_n)) \phi(x_n)^T$

$\Rightarrow \sum_{n=1}^N r_n t_n \phi(x_n)^T - \vec{w}^T \sum_{n=1}^N r_n \phi(x_n) \phi(x_n)^T = 0$

let  $t = [t_1, \dots, t_N]^T, \vec{\phi} = [\phi(x_1), \dots, \phi(x_N)]^T$

$R = (r_n, \dots, 0)^T \Rightarrow t^T R \vec{\phi} = \vec{w}^T \vec{\phi} (R \vec{\phi})^T$

$\Rightarrow w = (R^T R \vec{\phi})^{-1} R \vec{\phi}$

$\left[ \begin{array}{c} f(x) \\ f(x_n) \end{array} \right] \sim N\left( \begin{bmatrix} M(X) \\ m(X_n) \end{bmatrix}, \begin{bmatrix} k(x, x) & \cdots & k(x, x_n) \\ k(x_n, x) & \cdots & k(x_n, x_n) \end{bmatrix} \right)$  select  $M$  eigenvectors of SC cov matrix -- projection.

commonly use  $k(x_n, x_n) = E[f(x)f(x_n)] = \exp(-\frac{\beta}{2}\|x_n - x\|^2)$  covariance kernel

**GP Regression**: noise Gaussian noise  $p(t_n | f_n) = N(t_n | f_n, \beta^{-1})$  where  $\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$ ,  $S = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T$  constrain  $\|u\|=1$  (use Lagrange multiplier) to max

$t_n = f_n + \epsilon_n$  so given  $X = [x_1, \dots, x_N]$ ,  $t = [t_1, \dots, t_N]$ . the conditional is  $p(t|f) = N(t|f, \beta^{-1}I_N)$ .

**Marginal distribution GP**:  $p(f) = N(f|0, K)$   
 $p(t) = p(t|f)p(f) = N(t|0, C)$ , cov matrix  $C(x_n, x_m) = k(x_n, x_m) + \beta^{-1}\delta_{nm}$

**Covariance Function**:  $k(x_n, x_m) = \theta_0 \exp(-\frac{\beta}{2}\|x_n - x_m\|^2) + \theta_1 x_n^T x_m$

**Prediction**: predict  $t_{N+1}$  for new input  $x_{N+1}$   
 joint dist over  $t$  &  $t_{N+1}$  is  $P(t, t_{N+1}) = N(0, \begin{bmatrix} C_N & K \\ K^T & \beta C \end{bmatrix})$   
 More  $C_N$  is  $N \times N$  matrix with element  $C_N(x_n, x_m) = k(x_n, x_m) + \beta^{-1}\delta_{nm}$ ,  $c$  is scalar  $= k(x_{N+1}, x_{N+1}) + \beta^{-1}$   
 $K$  is  $N$  by 1 vector with elements  $k(x_n, x_{N+1})$   
 conditional distribution is Gaussian  $P(t_{N+1}|t) = N(m(x_{N+1}), \sigma^2(x_{N+1}))$  with  $m(x_{N+1}) = k^T C_N^{-1} t$ ,  $\sigma^2(x_{N+1}) = c - k^T C_N^{-1} k$

key results of GP:  
 hyperparameters ( $\beta$ )  
 max log of marginal likelihood wrt  $\theta$ :  $\ln p(t|\theta) = -\frac{1}{2} \ln |C_\theta| - \frac{1}{2} t^T C_\theta^{-1} t - \frac{N}{2} \ln(\pi)$   
 $\Rightarrow$  gradient wrt  $\theta$ :  $\frac{\partial \ln p(t|\theta)}{\partial \theta_i} = -\frac{1}{2} T_i (C_\theta^{-1} C_\theta) + \frac{1}{2} t^T C_\theta^{-1} C_\theta^{-1} t$

**Automatic Relevance Determination**:  $k(x_n, x_m) = \theta_0 \exp(-\frac{1}{2} \sum_{i=1}^D \gamma_i (x_{ni} - x_{mi})^2) + \theta_1 x_n^T x_m$ ,  $\ln p(t|\theta) = -\frac{1}{2} \ln |C_\theta| - \frac{1}{2} t^T C_\theta^{-1} t - \frac{N}{2} \ln(\pi)$

**Mixture models**:  $N \leftarrow \sum_{n=1}^N \Gamma_{nk} \|x_n - \mu_k\|^2$ , each data  $x_n$  as binary vector  $\Gamma_n$  of length  $K$  indicates which  $K$  cluster  $\mu_k$  is prototype. try to minimize J. (K-Means clustering): given  $\mu_k$ , (E-step)  $\text{rank}_k = \sum_{i=1}^N \Gamma_{ik}$  if  $k = \arg\min_j \|x_n - \mu_j\|^2$ ,  $p(z_k|x) = \prod_{i=1}^N \Gamma_{ik} N(x|\mu_k, \Sigma_k)$ , (M-step)  $\mu_k = \frac{\sum_{i=1}^N \Gamma_{ik} x_i}{\sum_{i=1}^N \Gamma_{ik}}$

**Joint dist.**:  $p(x, z) = p(x|z)p(z)$  sufficient stat: mean  $\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$ , sample cov.  $S$ . corresponding posterior  $p(z_k|x) = \prod_{i=1}^N \Gamma_{ik} N(x|\mu_k, \Sigma_k)$  whose obs are given by  $M$  principal eigenvectors of  $S$ .  $p(z_k|x) = \sum_{j=1}^M \Gamma_{kj} N(x|\mu_j, \Sigma_j)$   $L_m$  is  $M \times M$  diagonal containing  $M$  largest eigenvalues.  $R$  is arbitrary  $M \times M$  orthogonal matrix.

**ML for Gaussian mixture**:  $\ln p(X|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \sum_{k=1}^K \Gamma_{nk} N(x_n|\mu_k, \Sigma_k)$  EM for PCA:  $\log p(x, z|w, \Sigma) = \sum_{i=1}^N [\log p(x_i|z_i) + \log p(z_i)]$   $\Rightarrow$  diff.  $\mu_k = \frac{1}{\Gamma_{nk}} \sum_{i=1}^N \Gamma_{ik} x_i$ ,  $N_k = \sum_{i=1}^N \Gamma_{ik}$  zero noise limit: derive standard PCA as a limit of PCA  $\sigma^2 \rightarrow 0$

$\hat{\mu}_k = \frac{1}{N_k} \sum_{i=1}^N \Gamma_{ik} (x_i - \mu_k)^T$ ,  $\Gamma_{nk} = \frac{N_k}{N}$  ML para same. Posterior mean reduce to  $\lim_{\sigma^2 \rightarrow 0} (W^T W + \sigma^2 I)^{-1} W^T (x - \mu) = (W^T W)^{-1} W^T (x - \mu) \rightarrow$  standard PCA

**E: update  $\gamma(z_{nk})$** :  $\gamma(z_{nk}) \xrightarrow{\sigma^2 \rightarrow 0} \text{posterior covariance} = 0$ . Factor Analysis (FA): take PCA, joint & marginal are also Gaussian.

$p(z_{nk}=1|x) = \frac{\pi_k N(x|\mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x|\mu_j, \Sigma_j)}$  marginal:  $p(x) = \int_z p(z) p(x|z) dz = \sqrt{C/\pi} \int_{W^T z + \mu} W W^T + \Sigma dz$   $\xrightarrow{\text{max diag}}$

i: update  $\mu_k, \Sigma_k$ , don't forget check convergence & evaluate log-like.

**general EM**: give joint  $p(z, x|\theta)$ , max  $p(x|\theta)$  wrt  $\theta$ . E: compute posterior over latent  $p(z|x, \theta)$  add  $E(\theta) = \frac{1}{2} \sum_{i=1}^N \|\gamma(x_i, \theta) - \mu_i\|^2$

M: find new est of  $\theta$  new.  $\theta^{\text{new}} = \arg\max_{\theta} Q(\theta, \theta)$  Markov models: future predictions independent of all but the most recent observations.

where  $Q(\cdot) = \sum_{i=1}^N p(z|x, \theta^{\text{old}}) \ln p(x, z|\theta)$  1st order M&C  $p(x_1, \dots, x_N) = p(x_1) \prod_{n=2}^N p(x_n|x_{n-1})$  - conditional:  $p(x_n|x_1, \dots, x_{n-1}) = p(x_n|x_{n-1})$

**Gaussian Mixture Revist**:  $\ln p(x|\pi, \mu, \Sigma) = \prod_{n=1}^N \prod_{k=1}^K \Gamma_{nk} N(x_n|\mu_k, \Sigma_k)$  2nd order M&C  $p(x_1, \dots, x_N) = \prod_{n=1}^N p(x_n|x_1, \dots, x_{n-1})$

$\ln p(x, z|\pi, \mu, \Sigma) = \prod_{n=1}^N \prod_{k=1}^K \Gamma_{nk} N(x_n|\mu_k, \Sigma_k)$  State Space models (additional latent var introduced)

$\ln p(x, z|\pi, \mu, \Sigma) = \sum_{n=1}^N \sum_{k=1}^K \Gamma_{nk} \ln \Gamma_{nk} + \sum_{n=1}^N \ln N(x_n|\mu_k, \Sigma_k)$

$\ln p(x, z|\pi, \mu, \Sigma) \rightarrow \ln p(x|\pi, \mu, \Sigma) + \ln p(z|\pi, \mu, \Sigma)$   $\xrightarrow{\text{latent has M.C. if latent discrete (HMMs)}}$

$\ln p(x, z|\pi, \mu, \Sigma) \rightarrow \ln p(x|\pi, \mu, \Sigma) + \ln p(z|\pi, \mu, \Sigma)$  observed variables can be cont./discrete.

$\ln p(x, z|\pi, \mu, \Sigma) = \sum_{n=1}^N \sum_{k=1}^K \Gamma_{nk} \ln \Gamma_{nk} + \sum_{n=1}^N \ln N(x_n|\mu_k, \Sigma_k)$  joint  $p(x_1, \dots, x_N, z_1, \dots, z_N) = p(x_1) \prod_{n=2}^N p(z_n|z_{n-1}) \prod_{n=1}^N p(x_n|z_n)$

$\ln p(x, z|\pi, \mu, \Sigma) = \sum_{n=1}^N \sum_{k=1}^K \Gamma_{nk} [\ln \Gamma_{nk} + \ln N(x_n|\mu_k, \Sigma_k)]$  predictive dist:  $p(x_{n+1}|x_1, \dots, x_N)$  depends on all previous obs.

**HMM:** 1st order MC generates hidden state seq. (transition prob)  
 $p(z_n=k|z_{n-1}=j) = \pi_{jk}$ ,  $p(z_1=k) = \pi_{1k}$ ,  $\sum_{j=1}^K \pi_{jk} = 1$

conditionals:  $p(z_n|z_{n-1}, A) = \prod_{k=1}^K \prod_{j=1}^{z_{n-1}} \pi_{jk} \delta_{z_n=k}$ ,  $p(z_1|A) = \prod_{k=1}^K \pi_{1k}$

emission prob.  $p(x_n|z_n, \phi) = \prod_{k=1}^K \pi_{1k} p(x_n|\phi_k)$   $\xrightarrow{\text{(continuous)}}$   $\prod_{k=1}^K N(x_n|\mu_k, \Sigma_k)$  (mixture Gauss)