

SSE Loss Func. $E(\vec{w}) = \frac{1}{2} \sum_{n=1}^N (y(\vec{x}_n, \vec{w}) - t_n)^2 + \frac{\lambda}{2} \|\vec{w}\|^2$ reg. parameter. (don't minimize w!) \rightarrow decrease complexity, chose λ by cross validation

\rightarrow linear models \rightarrow linear in w . \rightarrow Least Squares $y(\vec{x}_n, \vec{w}) = \vec{x}_n^T \vec{w}$ (weight decays / Ridge reg.)

\rightarrow L_2 loss is $\vec{w}^T (\vec{x}_n \vec{x}_n^T) \vec{t}_n$, $x \in \mathbb{R}^{N+1}$ \rightarrow polynomial fit $y(\vec{x}, \vec{w}) = \sum_m w_m x_m^m$, M picked by Model Selection - Avoid overfitting, achieve generalization

\rightarrow assumed to be sampled iid so training data and test data have same dist!

\rightarrow Probabilistic (curve fit) $\vec{t} = y(\vec{x}, \vec{w}) + \vec{\epsilon}$, $\vec{\epsilon} \sim \mathcal{N}(0, \beta^{-1})$, $\vec{t} \sim \prod_{n=1}^N \mathcal{N}(t_n | y(\vec{x}_n, \vec{w}), \beta^{-1}) = p(\vec{t} | \vec{x}, \vec{w}, \beta)$, $\vec{w}_{ML} = E(\vec{w})$

\rightarrow $\vec{E}(\vec{w}) = \frac{1}{N} \sum_{n=1}^N (y(\vec{x}_n, \vec{w}) - t_n)^2$ \rightarrow minimize log likelihood w.r.t \vec{w} \rightarrow minimizing least squares.

RVs $Cov(\vec{x}, \vec{t}) = E\{[\vec{x} - E(\vec{x})][\vec{t}^T - E(\vec{t}^T)]\} \rightarrow \int \rho(t | \vec{x}, \vec{w}, \beta) dt$

$\rightarrow \int_{x \in [0, 1]} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1} dt$, $E(x) = \frac{a}{a+b}$, $\text{Var}(x) = \frac{ab}{(a+b)^2(a+b+1)}$, $\Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du$.

Beta Dist. $\text{Beta}(x | a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}$, $E(x) = \frac{a}{a+b}$, $\text{Var}(x) = \frac{ab}{(a+b)^2(a+b+1)}$, $\Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du$.

Multinomial $\text{Mult}(m_1, \dots, m_K | M, N) = \binom{N}{m_1, \dots, m_K} \prod_{k=1}^K M_k^{m_k}$ \rightarrow size of group, $\sum M_k = N$

Multivariate Gaussian $N(\vec{x} | \vec{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2}} \exp\left\{-\frac{1}{2}(\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})\right\}$, $\vec{x} \in \mathbb{R}^D$, $\vec{\mu} = \vec{U}(\vec{x} - \vec{\mu})$, $\Sigma_{kl} = \vec{x}_l \vec{x}_k$, so \vec{x} shifted by $\vec{\mu}$, rotated by \vec{U} , and stretched in U_l by $\sqrt{\lambda_l}$

$\hookrightarrow \ln p(\vec{x} | \vec{\mu}, \Sigma) = -\frac{ND}{2} \ln(\pi) - \frac{1}{2} \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N (\vec{x}_n - \vec{\mu})^T \Sigma^{-1} (\vec{x}_n - \vec{\mu})$, $\vec{w}_{ML} = \frac{1}{N} \sum_{n=1}^N \vec{x}_n$, $\Sigma_{ML} = \frac{1}{N} \sum_{n=1}^N (\vec{x}_n - \vec{\mu}) (\vec{x}_n - \vec{\mu})^T$ * of sufficient statistics

Mixture of Gaussians $p(\vec{x}) = \sum_{k=1}^K \pi_k N(\vec{x}; \vec{\mu}_k, \Sigma_k)$ Component

Exponential Family $\rightarrow p(\vec{x} | \vec{t}) = h(\vec{x}) g(\vec{t}) \exp(\vec{t}^T \vec{\phi}(\vec{x}))$, $h: \mathbb{R}^D \rightarrow \mathbb{R}$, $g: \mathbb{R}^D \rightarrow \mathbb{R}$, $\vec{\phi}: \mathbb{R}^D \rightarrow \mathbb{R}^M$, $\vec{t}: \mathbb{R}^D \rightarrow \mathbb{R}^M$, $-\nabla \log(g(\vec{t})) = E(\vec{v}(\vec{x}))$

Linear Basis Func. $\rightarrow y(\vec{x}, \vec{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\vec{x})$, $\phi_j: \mathbb{R}^D \rightarrow \mathbb{R}$ \rightarrow Can be nonlinear in input space \rightarrow Can make problem linearly separable \rightarrow Can be global or local.

Statistical Decision Th. $P(C_k | \vec{x}) = \frac{P(C_k) P(\vec{x} | C_k)}{P(\vec{x})}$ posterior \rightarrow prior \rightarrow likelihood \rightarrow marginal likelihood \rightarrow Can assign \vec{x} to C_k .

Loss Function $E[L] = \iint L(t, y(\vec{x})) p(\vec{x}, t) dx dt$, Minkowski Loss $\rightarrow E[L] = \iint (t - y(\vec{x}))^q p(\vec{x}, t) dx dt$, min of $q = \rightarrow$ Conditional median [LASSO] $q=2 \rightarrow$ E(t|x) [Squared Loss]

Bias-var tradeoff $E[L] = \iint (y(\vec{x}) - h(\vec{x}))^2 p(\vec{x}) dx + \iint (h(\vec{x}) - t)^2 p(\vec{x}, t) dx dt$, $h(\vec{x}) = E(t | \vec{x}) = \int t p(t | \vec{x}) dt$ bias \rightarrow Intrinsic variability of target (noise) variance \rightarrow Variance - Variability of model y to diff datasets D .

$\rightarrow E_D[y(\vec{x}, D) - h(\vec{x})]^2 = E_D[y(\vec{x}, D) - h(\vec{x})]^2 + E_D[y(\vec{x}, D) - E_D[y(\vec{x}, D)]]^2$ small λ \rightarrow high variance \rightarrow flexible model \rightarrow parameters can be big large λ \rightarrow high bias \rightarrow weight tied to D \rightarrow over regularized.

Bagging \rightarrow Compute via conjugate prior, MCMC... \rightarrow Use Gaussian approx \rightarrow Posterior distribution \rightarrow Likelihood \rightarrow Prior \rightarrow Marginal Likelihood \rightarrow Posterior \rightarrow Marginal Likelihood \rightarrow Bias \rightarrow Variance

Predictive Dist. $p(\vec{x} | D) = \int p(\vec{x} | D, \vec{w}) p(\vec{w} | D) d\vec{w} \rightarrow$ if $M=0$ \rightarrow Equivalent Kernel: $y(\vec{x}, \vec{m}_n)$ \rightarrow if $S_0 = \alpha^{-1} I$ \rightarrow $\vec{y}(\vec{x}) = \sum_{n=1}^N K(\vec{x}, \vec{x}_n) t_n$ \rightarrow $t_n = \vec{x}_n^T \vec{\phi}(\vec{x})$ \rightarrow uncertainty of parameter vals. \rightarrow $N \rightarrow S_0 \rightarrow$ $S_0 \rightarrow$ $\vec{w} \rightarrow$ \vec{w}_{post}

Bayesian Linear Regression $P(\vec{w} | \vec{t}, \vec{x}, \beta) = \prod_{n=1}^N N(t_n | \vec{w}^T \vec{\phi}(\vec{x}_n), \beta^{-1}) N(\vec{w} | \vec{m}_0, S_0)$ Posterior over weights \rightarrow Marginal Likelihood \rightarrow Conjugate prior \rightarrow Average predictive dist. weighted by its posterior \rightarrow This is a mixed model. \rightarrow Model Selection - use model with most model evidence $p(D | M_1) \approx p(D | M_2)$

Bayesian Model Comparison $P(M_i | D) \propto P(M_i) p(D | M_i)$, $p(t | x, D) = \sum_i p(t | x, M_i, D) p(M_i | D)$ \rightarrow $\ln p(D, M_i) \equiv \ln p(D | w_{map}) + \ln \left(\frac{\Delta w_{post}}{\Delta w_{prior}} \right)$ \rightarrow maximize marginal likelihood to determine the value of the hyperparameters from the training data.

Evidence Approx. $P(t | \vec{x}, D) = \iint P(t | \vec{x}, \vec{w}, \beta) P(w | D, \alpha, \beta) P(\alpha, \beta | D) dw d\beta$ precision of output noise precision of prior training data \rightarrow secret (α, β) \rightarrow MLE of $P(t | \vec{x}, \vec{w}, \beta) \propto P(\alpha, \beta | D)$ [assuming flat prior $P(\alpha, \beta)$]

Classification **Linear** $\text{① Discriminant function: } y_K(\vec{x}_n) = \vec{x}_n^T \vec{w}_K + w_0 = \vec{w}^T \vec{x}_n$ \rightarrow K classes, $\vec{x} = \begin{bmatrix} \vec{x}_1 \\ \vdots \\ \vec{x}_N \end{bmatrix} \in \mathbb{R}^{N \times (d+1)}$ \rightarrow $\vec{w}^T \vec{x}_n$ is one of K encoded \rightarrow optimal $(\vec{x}^T \vec{x})^{-1} \vec{x}^T \vec{t}$ $\in \mathbb{R}^{(d+1) \times K} \rightarrow$ Assign \vec{x}_n to $\max_K y_K(\vec{x}) = \max_K \vec{w}^T \vec{x}_n$ \rightarrow WXT \rightarrow Sensitive to outliers, bad w.r.t multiple classes. \rightarrow Maximize b/w class variance \rightarrow want large b/w means and little spread.

Fisher's Linear Discriminant - project onto vector $(M_2 - M_1)^\perp$ \rightarrow minimize within-class variance \rightarrow want large b/w means and little spread. \rightarrow These directly map input into a class. \rightarrow minimize within-class variance \rightarrow want large b/w means and little spread. \rightarrow Softmax \rightarrow we assume $\sim N_0(M_K, \Sigma^{-1})$ then linear in input + prior \rightarrow $\text{softmax} = \frac{p(\vec{x} | c_k) p(c_k)}{\sum_i p(\vec{x} | c_i) p(c_i)} = \frac{\exp(\alpha_k)}{\sum_i \exp(\alpha_i)}$ $\alpha_k = \ln(p(\vec{x} | c_k) p(c_k))$ space!

② Generative Approach - model $p(\vec{c} | \vec{x})$ indirectly by modelling $p(\vec{x} | \vec{c})$ \rightarrow ML Soln - $P(\vec{t}, \vec{x}, \vec{t}, \vec{c}, \vec{f}, \vec{t}_2, \vec{s}) = \prod_{n=1}^N [\pi N(\vec{x}_n | M_1, \Sigma_1)]^{1-t_n} [\pi N(\vec{x}_n | M_2, \Sigma_2)]^{1-t_n}$, $\Pi_{ML} = \frac{M_1}{N_1 + N_2}$, $M_{i, ML} = \frac{1}{N_i} \sum_{n=1}^N t_n \vec{x}_n$, $S_{ML} = \frac{1}{N} \sum_{n=1}^N (\vec{x}_n - M_1)(\vec{x}_n - M_1)^T$

$S_{ML} = \frac{N_1}{N} S_1 + \frac{N_2}{N} S_2 \rightarrow$ Note: we have to assume data comes from Normal dist.!

③ Discriminative Approach - model $P(C_k | \vec{x})$ directly. OR fix basis functions \rightarrow Functions, $\phi(\vec{x}) \rightarrow$ Decision boundary linear in feature space ϕ but not in input space. \rightarrow logistic Regression - $P(C_k | \vec{x}_n) = \frac{\exp(\vec{w}_k^T \vec{x}_n)}{1 + \exp(\vec{w}_1^T \vec{x}_n)}$ \rightarrow Likelihood - $P(T | \vec{x}, \vec{w}_1, \dots, \vec{w}_K) = \prod_{n=1}^N \left[\prod_{k=1}^K p(C_k | \vec{x}_n)^{t_{nk}} \right] = \prod_{n=1}^N \left[\prod_{k=1}^K \frac{1}{1 + \exp(-\vec{w}_k^T \vec{x}_n)} \right]$

\rightarrow -ve log likelihood = C.E. = $-\sum_{n=1}^N (y_{n1} - t_{n1}) \vec{x}_n^T$ (for \vec{w}_1) \rightarrow (2 class) $\logistic \text{ Regression} = \prod_{n=1}^N y_n t_n (1 - y_n)^{1-t_n} \frac{\partial}{\partial w_1} = y_n (1 - y_n) \vec{x}_n$ since $\frac{\partial}{\partial a} \sigma(a) = \sigma(a)(1 - \sigma(a))$

$\sigma(\vec{w}^T \vec{x}) = \frac{1}{1 + \exp(-\vec{w}^T \vec{x})}$, $\sigma(-a) = 1 - \sigma(a)$, $\sigma^{-1}(a) = \ln(\frac{a}{1-a})$ $\left[\begin{array}{c} C \\ \vec{w} \\ \vec{r} \\ c \end{array} \right]$

GP $\vec{f} = \begin{pmatrix} f(x_1) \\ \vdots \\ f(x_N) \end{pmatrix} \quad [X]_{nm} = K(x_n, x_m), \quad \vec{f} \sim N(\vec{0}, K), \quad \vec{f} \sim N(\vec{f}, \beta^{-1} I), \quad \vec{f} = \int \vec{t} (\vec{t}^T \vec{f}) p(\vec{f}) d\vec{f} \sim N(\vec{0}, \Sigma), \quad \left[\begin{array}{c} \vec{t} \\ t_{N+1} \end{array} \right] \sim N(\vec{0}, \Sigma), \quad \left[\begin{array}{c} K(x_n, x_1) & \dots & K(x_n, x_{N+1}) \\ K(x_1, x_1) & \dots & K(x_1, x_{N+1}) \\ \vdots & \ddots & \vdots \end{array} \right]$

$[C]_{nm} = K(x_n, x_m) + \delta_{n,m} \beta^{-1}$ \rightarrow Marginal of f \rightarrow Conditional on f \rightarrow Marginal of \vec{t} , using that $\vec{f} \sim \vec{t}$ \rightarrow predict t_{N+1} from \vec{x} and x_{N+1} and \vec{t} gaussian!

$t_n = f_n + \epsilon_n$, $\epsilon_n \perp \epsilon_i$, $\epsilon_n \sim N(0, \sigma^2)$ \rightarrow inverting this is $O(N^3)$ \rightarrow data points or basis functions! \rightarrow And you get Bayesian Linear Regression!

finally, $t_{N+1} | \vec{t} \sim N(\vec{t}^T \vec{C}^{-1} \vec{t}, \sigma^2 \vec{C}^{-1} \vec{t})$, if $\vec{f} = \Phi \vec{w}$, $\vec{w} \sim N(\vec{0}, \alpha^{-1} I) \Rightarrow [K]_{nm} = \frac{1}{\alpha} \phi(\vec{x}_n)^T \phi(\vec{x}_m)$, $f \sim N(\vec{0}, K)$

\rightarrow $\sum_{n=1}^N \left[\sum_{k=1}^N K(x_n, x_k) t_k \right] \leftarrow$ linear combo \rightarrow can use parametric family of Covariance functions and infer of Targets parameters from data \rightarrow learning hyperparameters using Type II ML in standard model!

Hyper Parameters

Type I MLF: Maximize $p(\theta|t) = \frac{p(t|f)p(f)^{\text{prior}}}{p(t)}$ over params, f . Type II: $\max_{\theta} p(t|\theta) = \int p(t|f)p(f|\theta) df$, integrate over params, maximize over θ .

(Final) introduce $P(\theta)$, infer $p(\theta|t)$ posterior using MCMC, since not tractable!

K-Means $\vec{x}_1, \dots, \vec{x}_N \in \mathbb{R}^D$ into K clusters, minimizes $\sum_{n=1}^{NK} \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$. Given M_K , set $r_{nk}=1$ if $k = \arg\min_j \|x_n - \mu_j\|^2$, $r_{nk}=0$ otherwise [put x_n in cluster M_K]. Set M_K to mean of x_n 's in M_K .

LATENT VARIABLE MODELS

- Continuous \rightarrow PCA
- Discrete \rightarrow MoG
- Mixture Models \rightarrow HMM

HMM: Data $\vec{x} \in \mathbb{R}^{NxD}$ (N latent vars, 1 for each x_n)

EM Algorithm: Given x , $p(z|x, \theta)$, $p(x|z|\theta)$ [Goal] find θ s.t. $p(x|\theta)$ maximized! Notice $\ln p(x|\theta) = \ln \left(\sum_z p(x|z|\theta) \right) \rightarrow$ Marginal log likelihood \rightarrow bad

E: Compute $E_z[\ln p(x|z|\theta)] = \sum_z \ln p(x|z|\theta) p(z|\theta)$ [old] Express value of $p(z|\theta)$ given complete data LL.

M: Find $\max_{\theta} Q(\theta, \theta^{old})$, set θ^{old} to it. (for Gaussian case) $Q(\theta, \theta^{old}) = E_z \left\{ \sum_{n=1}^N z_{nk} [\ln \pi_{nk} + \ln N(x_n | \mu_k, \Sigma_k)] \right\} = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) [\ln \pi_{nk} + \ln N(x_n | \mu_k, \Sigma_k)]$

Computations: $E(z_{nk}) = \sum_{n=1}^N \gamma(z_{nk}) p(z_{nk} = 1 | \vec{x}_n) = \gamma(z_{nk}) / M_j^{new} = \frac{2}{2\pi} \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) [\ln \pi_{nk} + \ln N(x_n | \mu_k, \Sigma_k)] = \sum_{n=1}^N \gamma(z_{nk}) (\vec{x}_n - \vec{\mu}_j)^T \Sigma_j^{-1} \vec{x}_n$ [multiply both sides by $\vec{\Sigma}_j^{-1}$]

$$\Rightarrow \sum_{n=1}^N \gamma(z_{nk}) \vec{x}_n - \sum_{n=1}^N \gamma(z_{nk}) \vec{\mu}_j = M_j^{new} \Rightarrow \frac{\sum_{n=1}^N \gamma(z_{nk}) \vec{x}_n}{\sum_{n=1}^N \gamma(z_{nk})} = M_j^{new} = \frac{2}{2\pi} \sum_{n=1}^N \left[\gamma(z_{nk}) \ln(\pi_{nk}) + \sum_{k=1}^K \gamma(z_{nk}) \ln(1 - \pi_{nk}) \right] = \sum_{n=1}^N \frac{\gamma(z_{nk})}{\pi_{nk}} - \sum_{k=1}^K \frac{\gamma(z_{nk})}{1 - \pi_{nk}} = 0 \Rightarrow \frac{N_j}{\pi_{nj}} = \frac{N - N_j}{1 - \pi_{nj}} \Rightarrow \pi_j^{new} = \frac{N_j}{N}$$

EM = $E_z[\ln p(x|z|\pi, M, \Sigma)] \rightarrow -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \|x_n - \mu_k\|^2 + C = \text{K-means}$ [When $\epsilon \rightarrow 0$, $\gamma = \mathbb{I}$ shared]

PCA: Select top eigenvectors of $S = \{u_1, \dots, u_m\}$, $z_n = x_n^T u_i$ [project data] ① Max Var Formulation - let $x \in \mathbb{R}^{NxD}$ $u_i^T u_i = 1$, want to maximize $\frac{1}{N} \sum_{n=1}^N (u_i^T x_n - u_i^T \bar{x})^2 = u_i^T S u_i$ wrt u_i , constraint $\|u_i\|=1 = \frac{\partial}{\partial u_i} u_i^T S u_i + \lambda(1 - u_i^T u_i) = \gamma_i u_i \Rightarrow \gamma_i$ eigenval of S . ② Min Error Formulation - use orthogonal basis $\{v_1, \dots, v_D\}$ $\tilde{x}_n = \sum_{i=1}^m z_n u_i + \sum_{i=m+1}^D b_i v_i$

$x_n = \sum_{i=1}^D (x_n^T v_i) v_i$ want to minimize $\frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2$ to get $z_n^{min} = x_n^T v_i \Leftrightarrow \sum_{i=1}^D v_i^T S v_i$ min this $\Leftrightarrow \sum_{i=1}^D \lambda_i$ [minimized by choosing smallest λ_i e-vals].

PCA for DT: When $N < D$, from S reduced to $O(N^3)$, when X centered $\frac{1}{N} X^T X u_i = \lambda_i u_i \Rightarrow \frac{1}{N} X X^T (X u_i) = \lambda_i X u_i \Rightarrow X u_i$ e-val for $X X^T$ w/ same e-val λ_i !

PPCA: Generative, assume $p(z|x) = N_0(x | Wz + m, \sigma^2 I)$, density/prob. dist. $\rightarrow p(x) = \int p(z)p(x|z) dz = N(x | m, W W^T + \sigma^2 I)$ $E(x) = E(m + Wz + \epsilon) = m$, $C(x) = E((m + Wz + \epsilon)^2 - m)(Wz + \epsilon)^T$

* rotations in latent space $p([z]) = N([z] | [m] [I_w \vec{w}^T \sigma^2 I])$, $p(z|x) = N(M^{-1} W^T (x - m), \sigma^2 M^{-1})$ $M = W^T W + \sigma^2 I$ \rightarrow can now solve for W and σ^2 directly with ML or use EM algorithm!

EM PCA: faster than PPCA, use log likelihood $\ln(p(x|z|m, w, \sigma^2)) = \sum_n [\ln p(x_n|z_n) + \ln p(z_n)]$ * Does not require constructing cov matrix. $| \sigma^2 \rightarrow \text{ppca} \rightarrow \text{PCA}$ ML params same!

FA: $p(x|z) = N(wz + m, \psi)$, $\psi \in \mathbb{R}^{DxD}$ diagonal, $p(x) = N(m, W W^T + \psi)$

Factors nonidentifiable, can find many sets of params to get some ML \rightarrow degeneracy interpretation of factors impossible!

Markov model w/ random members

Name	Solve directly w/ ML?	uses for directions of large rotation?	scale invariant?
PCA	Yes	Variance	Yes
FA	No	Correlations	No Yes

Auto-encoder: Neural nets w/ same # of inputs as outputs. $M < D$ hidden units (bottleneck) \rightarrow nonlinear extension of PCA.

HMM: x_i is not iid $\rightarrow p(x_n | x_1, \dots, x_{n-1}) = p(x_n | x_{n-1})$ [1st order] $(z_1) \rightarrow (z_2) \rightarrow \dots \rightarrow (z_{n-1}) \rightarrow (z_n)$ mixture model w/ states completed over time $A = \begin{bmatrix} P(\text{from 1 to 1}) & P(\text{from 1 to 2}) \\ P(\text{from 2 to 1}) & P(\text{from 2 to 2}) \end{bmatrix}$ $\xrightarrow{\text{all obs same}} \text{Standard mixture model!}$

Markov Assumption: $p(x_1, \dots, x_n) = p(x_1) \prod_{i=2}^n p(x_i | x_{i-1})$

State Space Model: introduce latent vars which form Markov chain. $\Pi_{jk} = p(z_{ik} = 1) \rightarrow [A]_{jk} = p(z_{ik} = 1 | z_{i-1} = 1), \sum_k A_{ik} = 1$

Using $z_{n+1} \perp z_{n-1} | z_n$

EM for HMM's: note $p(x|\theta) = \sum_{z_{1:n}} p(x|z|\theta) = \sum_{z_1} p(z_1) p(x_1|z_1) \dots \sum_{z_n} p(z_n|z_{n-1}) p(x_n|z_n)$, usually $p(x|\theta) = \sum_{z_1} p(z_1) p(x|z_1)$ since z_1 indep! joint prob of data up to point $z_n + z_{n+1}$

E: Compute $\gamma(z_n) = p(z_n|x, \theta^{old}) = \frac{p(x|z_n)p(z_n)}{p(x)} = \frac{p(x_1, \dots, x_n|z_n)p(z_n)p(x_{n+1}, \dots, x_N|z_n)}{p(x)} = \frac{p(x_1, \dots, x_n, z_n)p(x_{n+1}, \dots, x_N|z_n)}{p(x)} = \frac{\sum_{z_{n+1}} p(z_{n+1}|z_n) \beta(z_n)}{p(x)}$ after z_n .

$\alpha(z_n) = p(x_n|z_n)p(x_1, \dots, x_{n-1}, z_n) = p(x_n|z_n) \sum_{z_{n-1}} p(x_{n-1}, \dots, x_1, z_{n-1}, z_n) = p(x_n|z_n) \sum_{z_{n-1}} p(x_{n-1}, \dots, x_1, z_{n-1}) p(z_n|z_{n-1}) = p(x_n|z_n) \sum_{z_{n-1}} \alpha(z_{n-1}) p(z_n|z_{n-1})$

$\beta(z_n) = \sum_{z_{n+1}} p(x_{n+1}, \dots, x_N, z_{n+1}, z_n) p(z_{n+1}|z_n) = \sum_{z_{n+1}} p(x_{n+1}, \dots, x_N|z_{n+1}) p(z_{n+1}, z_n) = \sum_{z_{n+1}} p(x_{n+1}, \dots, x_N|z_{n+1}) p(x_{n+1}|z_{n+1}) p(z_{n+1}|z_n) = \sum_{z_{n+1}} \beta(z_{n+1}) p(x_{n+1}|z_{n+1})$

also $\alpha(z_i) = p(x_i, z_i) = p(z_i) p(x_i|z_i) = \prod_{k=1}^K \pi_k p(x_i | \phi_k)$, $\beta(z_n) = p(z_n|x) = \frac{p(x, z_n|\theta)}{p(x)} = \frac{p(x, z_n|\theta)}{\sum_{z_n} p(x, z_n|\theta)} \rightarrow \beta(z_n) = 1$

Notice $\sum_{z_n} \frac{\alpha(z_n)\beta(z_n)}{p(x)} = \sum_{z_n} p(z_n|x, \theta) = 1 \rightarrow p(x|\theta) = \sum_{z_n} \alpha(z_n)\beta(z_n) = \sum_{z_n} \alpha(z_n)$ $\xrightarrow{\text{set } n=N, \beta(z_N)=1}$

$E(z_{n-1}, z_n) = p(z_{n-1}, z_n|x) = p(x|z_{n-1}, z_n) p(z_{n-1}, z_n) = \alpha(z_{n-1}) p(x_n|z_n) p(z_n|z_{n-1}) \beta(z_n)$

M: ② maximize expected complete data log likelihood $\sum_z p(z|x, \theta^{old}) \log p(x|z|\theta)$ for θ and replace θ^{old} w/ θ .

$Q(\theta, \theta^{old}) = \sum_{k=1}^K \gamma(z_{ik}) \log \pi_k + \sum_{n=2}^N \sum_{k=1}^K \sum_{j=1}^K \gamma(z_{ik}) \log A_{jk} + \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{ik}) \log p(x_n|z_k)$, $\pi_k^{new} = \frac{\gamma(z_{ik})}{\sum_j \gamma(z_{ij})}$

Viterbi Decoding: choose $\gamma(z_n)$ w/ largest probability instead of sums to get path of max expected # of correct states.

HMM Ext.: Autoregressive \rightarrow better at catching long range correlations

- Input-Output HMM \rightarrow emission probs and transition probs depend on sequence u_1, \dots, u_N
- Factorial HMM \rightarrow 2 Markov chains of latent variables \rightarrow Exact inference intractable
- Regularizing HMMs: high dim outputs = lots of params in emission model \rightarrow \uparrow training data, tie params across states.
- Many States = transition matrix has many entries \rightarrow force some transitions to $\equiv 0$

prior - I want all weights = 0
likelihood - I want to fit data!