Machine Learning study sheet (dense handwritten notes). Selected legible content transcribed below; much of the page is extremely dense handwriting.

**Bayes rule:** $P(M_i|x, f_i) = \dfrac{P(f_i|x, M_i) \, P(x|M_i)}{P(f_i|x)}$

**Categorical Naive Bayes classifier:** $x$ discrete too; $P(x_c = i | Y = c) = \Theta_{ci}$, MLE.

**Linear classifier:** $y = \text{sign}(w^T f(x))$, a non-linear $y = \text{sign}(w^T f(x))$.

**Neural Networks:** linear unit $\varphi(x) = x$; sigmoid/logistic $\sigma(z) = \dfrac{1}{1+e^{-z}}$; $G(z)$; tanh $\tanh(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$; ReLU.

**Maximum likelihood:** $\arg\max_\theta \sum \log P(y_i|x_i, \theta)$.

**MAP:** $\arg\min -\log P(y|x)$, with prior, regularization.

**Bayesian decision theory.**

**Confusion matrix.** Precision: $\dfrac{TP}{TP+FP}$, Recall: $\dfrac{TP}{TP+FN}$. F-score $= \dfrac{2 \cdot TP}{2TP + FP + FN}$, harmonic mean. Accuracy $= \dfrac{TP+TN}{TP+TN+FP+FN}$.

**ROC curve** (Recall/TPR vs FPR). **Precision-Recall curve.**

**Gradient Descent.**

**Regularization:** ridge regression, LASSO.

**Perceptron:** $w' = w + y_i x_i$.

**SGD** (Stochastic Gradient Descent).

**Support Vector Machine (SVM):** hinge loss, margin. Kernels: Gaussian, Laplacian, polynomial.

**k-Nearest Neighbor.**

**Logistic Regression.**

**Naive Bayes classifier.**

**Gaussian Naive Bayes (GNB), General Gaussian Naive Bayes (GGNB).**

**LDA / QDA** (Linear / Quadratic Discriminant Analysis).

**Markov chains, Transition matrix, Stationary distribution.**

**Hoeffding's inequality, Confidence bounds, Cross-validation.**

**Jacobians, Backpropagation.**

**Sequence to sequence: HMM (Hidden states)** ⟹ generative; $X_{1:T}$ observations, $Y_{1:T}$ hidden states; $X,Y$ categorical. Same $p_y$ and $\theta_{y,x}$, additional $P(X_t = x_t | Y_t = y_t)$.

MLE: $\hat{p}_{X_1,Y_1} = \frac{\#[X_1=x, Y_1=y]}{\#[Y_1=y]}$

Inference tasks: ← Filtering: $P(Y_t | x_{1:t})$. Prediction: $P(X_{t+1} | x_{1:t}) = \sum_y P(X_{t+1} = y | x_{1:t})$. Conditioning: $P(Y_{t+1} | x_{1:t+1}) \propto P(Y_{t+1} | x_{1:t})$. Smoothing: $P(Y_t | x_{1:T}) \to \frac{1}{\sum_y} P(x_{t+1:T} | Y_t) P(Y_t | x_{1:t})$.

2. Prediction: $P(X_{t+1} | x_{1:t})$. Prediction: $P(X_{t+1} | x_{1:t}) = \sum_y P(X_{t+1} = y | Y_t = y) P(Y_t = y | x_{1:t})$

Matrix/Vector notation of filtering: $P(Y_t | x_{1:t})$ as vector $p^{(t)}_y P(Y_t = y | x_{1:t}) \Rightarrow p^{(t)}$ as vector of $P(Y_t = y | x_{1:t})$. Transition matrix $T$ as before, $\theta = [P(x_t | y_t)...P(x_t | y_t)]$; vector repr. of $P(X_t = x_t | Y_t = y)$

**Markov property**

For $t=1:T$: $P(Y_t = y | x_{1:t})$ as vector $p^{(t)}$: $\frac{P(Y_t = y | x_{1:t}) \cdot P(Y_t = y | x_{1:t})}{\alpha(y)}$

**Smoothing:** $P(Y_t = y | X_{1:T}) = \sum_y P(X_{t+1:T} | Y_t = y, X_t) \cdot \frac{P(Y_t = y | x_{1:t})}{b(y)}$

**Forward pass:** Init: $\alpha_1(y) = p_y \cdot \theta_{x_1,y}$. For $t=2:T$: $\alpha_t(y) = \sum_{y'} \alpha_{t-1}(y') \theta_{y',y} \cdot \theta_{x_t,y}$

**Backward pass:** Init: $b_T(y) = 1 \forall y$. For $t=T-1:1$: $b_t(y) = \sum_{y'} b_{t+1}(y') \theta_{y',y'} \cdot \theta_{x_{t+1},y'}$

**Forward-Backward - Algorithm:** $O(T \cdot c^2)$

For $t=1:T$: $P(Y_t = y | x_{1:T}) = \frac{1}{z} \alpha_t(y) b_t(y)$ where $z = \sum_y \alpha_t(y) b_t(y)$

**Active Learning:** expensive labels (need to ask expert)! Uncertainty sampling (p=0.5): $\hat{x}$ estimate at $\beta(y|x)$ given current labeled data. 2. pick most uncertain ($p\approx0.5$). 3. expect labels $y_i$. → till assumption violated. **Unsupervised Learning:** (anomaly)=doesn't fit any cluster, clustering = classification, reduction=unsupervised regression.

**Dim Reduction**

**Principal Component Analysis (PCA):** linear fixed $x$, the dimension reduction of $\to \lambda$. A. 1: $\max_z \|x z + A_z z^T x_z - x_z\|^2$, $A_z z^T = I_k$, $z^* = \arg\min_z \|x_z A\| \to x_i$ reconstruction=low dim approx.

$z = \frac{1}{n} \sum_i x_i^T x_i x_z$, assume centered: $z = \sum_i x_i = 0$, reconstruction loss minimization: dimension reduction. $(w_i^*, z_i^*) = \arg\min_{w,z} \sum_i \|w z_i - x_i\|^2$, $z_i = w^T x_i$, now let $x \to \phi(x) \Rightarrow w = \sum_i \alpha_i \phi(x_i)$ then normal PCA gets $\max \sum_i (\frac{z_i}{\|w\|}\phi(x_i))^2$

Variance maximization: $W^* = \arg\max_W W^T \Sigma W$, $W^T W = I_k$. Solve via SVD: $\Sigma = \frac{1}{n} A^T A = V \Lambda V^T$, $\Lambda_1 \geq ... \geq \Lambda_d$. Then $W^* = (V_1, ..., V_k) = $ first $k$ eig columns of $V$ corresponding to top $k$ eigenvalues of $\Sigma$ = principal component and then $z_i = W^T x_i$. $k$ chosen visually or as in $k$-means heuristics.

**Variance via max mean:** like SVD: $\Sigma = \sum_i \lambda_i V_i V_i^T$, $A = V \Sigma V^T = U \Sigma V^T$, $U = V \Sigma U^T$ ($\Sigma = V^T U^T$). $W^* = \arg\max_W W^T \Sigma W$

**Nonlinear methods:** kernel PCA: for kernels we require $w = \sum_i \alpha_i x_i$. From normal PCA: $\Lambda = $ eigenvalue of $X^T X$. $X^T X w = \Lambda w$, define $X w = \frac{1}{\lambda} w \Rightarrow w = \frac{X^T w}{\lambda} = \sum_i \frac{\alpha_i}{\lambda} x_i$, now let $x \to \phi(x) \Rightarrow w = \sum_i \alpha_i \phi(x_i)$, then normal PCA gets $\max \sum_i (\frac{z_i}{\|w\|}\phi(x_i))^2$

**Neural Network Autoencoders:** try to learn the identity fct: $d$ inputs, $d$ output units. $(W, W') = \arg\min_{W,W'} \frac{1}{n} \sum_i \|W' \phi(W x_i) - x_i\|^2$ one layer with $kd$ units. Initialization matters and is challenging. $P(x_i | z) \to PCA$ $\Rightarrow$ SGD backpropagation.

**Locally Linear Embedding (LLE):** A. identify neighbors $\theta N(i)$ of $i$. 2. weighted comb of neighbors: $x_i^* \approx \sum_{j \in N(i)} w_{i,j} x_j$ where $W = \arg\min_W \frac{1}{n} \sum_i \|x_i - \sum_{j \in N(i)} w_{i,j} x_j\|^2$, s.t. $\sum_i w_{i,j} = \lambda$ and $w_{i,j} = 0$ if $x_j \notin N(i)$. 3. Project to low-dim space: $z_{1:n} = \arg\min_z \sum_i \|z_i - \sum_{j \in N(i)} w_{i,j} z_j\|^2$, s.t. $\sum_i z_i z_i^T = I_k$, $\sum_i z_i = 0$, convex.

**Multi-dimensional scaling (MDS):** Given: dissimilarity matrix $D$ (i.e. city distances): $n \times n$. Find $z_{1:n}$ s.t. $\|z_i - z_j\| \approx D_{i,j}$ in 2D (=map) non-convex! (iterate map=equivalent=local opt)

$z_{1:n} = \arg\min_z \sum_{i,j} w_{i,j} (\|z_i - z_j\| - D_{i,j})^2$, where $w_{i,j} = \lambda$ =unnormalized, or $w_{i,j} = \frac{1}{D_{i,j}}$

**Center the kernel:** $E = \frac{1}{n} \mathbb{1} \mathbb{1}^T$

**Clustering: k-means clustering:** $\arg\min_{\mu_1...\mu_k} \sum_{i} \min_{j} \|x_i - \mu_j\|^2$. $L_{(\mu,z)} = \frac{1}{n} \sum_i \|x_i - \mu_{z_i}\|^2$, $\mu^* = \arg\min_\mu L_{(\mu,z)}$, non-convex. Lloyd's heuristic: Init $\mu^{(0)}_{1:k}$. convex

$\Rightarrow f(x,y) \propto \exp(-\frac{A}{2}(x-\mu_x)^2)$, $f_{Y|X}(x) = \frac{A}{\sqrt{2\pi G^2}} \exp(-\frac{(y-c)^2}{2G^2})$

**Gaussian math:** $f(x,y) \propto \exp(-\frac{A}{2 G^2}(x-\mu_x)^2)$, $f_{Y|X}(x) = \frac{A}{\sqrt{2\pi}} \exp(-\frac{(y-c)^2}{2 G^2})$, $a = \alpha(x), b,c,d = f$

**Assignment step:** $z_i^{(t)} = \arg\min_j \|x_i - \mu_j^{(t-1)}\|^2$ assigns to nearest cluster. $\Sigma_j = \frac{1}{n_j} \sum_{i:z_i=j} (x_i-\mu_j)(x_i-\mu_j)^T$

**Update step:** $\mu_j^{(t)} = \frac{1}{n_j} \sum_{i:z_i=j} x_i$ → mean of assigned points. Complexity $O(\text{iter}\cdot n \cdot k \cdot d)$

**A. 2:** $z_i = \arg\min_j \|x_i - \mu_j\|^2$. $(\text{M})z_{1:n} = \arg\min_{z} \sum_i \|x_i - \mu_{z_i}\|^2 \Rightarrow$ PCA with diff constraints

A. z: $z_i = \arg\min_j \|x_i - \mu_j\|^2$; the 2-D fine Cost O/C etc.

$\Rightarrow f(x,y) \propto \exp(\frac{1}{2}(x-\mu_x)^2) \cdot \sum_j f_{Y|X}(y|0) \propto f_{Y|X}(x,y) \cdot f(\frac{x_1}{x_2}) \propto$ rearrange to Gaussian

**Multivariate:** $f(x;\mu) = (2\pi)^{-d/2} |\Sigma|^{-1/2} \exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu))$

$\Rightarrow X = 2X_1 = 2X_2$. $f_{Y|X}(x;0) \propto f_{Y|X}(x,y) \cdot f(\frac{x_1}{x_2}) \propto$ same $\Sigma$ var, $\Rightarrow$ same $\Sigma$

$\propto \exp(-\frac{\lambda}{2}|x-\mu|) = \kappa_{1,k}\Sigma$ given $Y = X_1 + X_2$, $\kappa_3 = 0$

**Mixture Models:** Bayes classifier with unknown labels, generative $\Rightarrow$ can detect outliers, generate new point, sample cluster $y_c \sim$ categorical with indices $\Rightarrow$ equivalent solution $\Rightarrow$ multiple local optima! also for GMM: $P(X|Y=y_j) = w_j$. Then sample $x_i \sim f(x;\theta_j)$. non-convex (exchange components), hard to maintain constraint of symm. pod. variances, parameters $\Rightarrow$ non-spherical shapes. $k, U$ via CV (works here) or as in $k$-means. GMM degeneracy: $m = $ # components: if $m > n$, place each $z_j$ on a data point $x_j$, overfitting: $n z_j$: log-likelihood $\to \infty$ when var and $\to 0 \Rightarrow$ either reset any component that falls on a data point or use kickstart pair for add $\alpha I$ s.e.s $\approx MH$ to the cov an e matrices. Bayes

**Hard EM:** Compute hard assignments $y_i = \arg\max_y P(y|x_i, \theta^{(t)}) = \arg\max_y P(y|\theta^{(t)}) P(x_i|y, \theta^{(t)})$

**M-Step:** rough est. labels $\Rightarrow$ Bayes classifier: $\theta^{(t)} = \arg\max_\theta P(x_{1:n}, x_{1:n}|\theta)$. Hard EM bad, bc: if components overlap, and also: model is uncertain, but output are fixed!

**EM for GMM** with hard assignments $y_i = \arg\max_j f(x_i, \theta^{(t)})$, $\sigma^2 = $ variance $\sigma^2 \to 0$, uniform weights $w_j = \frac{1}{m}$

and identical spherical covariances $\Sigma_j = \sigma^2 I \Rightarrow k$-means.

**Soft EM:** Introduce latent variable $Y_i \in \{a,...,m\}$ (=m components), $P(Y=j|x_i) = w_j$, $w_{1:m} \geq 0$, $\sum w_j = 1$

**MLE:** $\forall$ Mixture Models: $w_j^* f(x_i; \theta_j)$, $GMM: f(x_i; \theta_j) = N(x_i; \mu_j, \Sigma_j)$

**MLE:** likelihood for one data point: $P(x_i|\theta) = \prod_j w_j f(x_i; \theta_j)$

Complete likelihood $P(X|\theta) = \prod_i \sum_j P(Y=j|\theta) P(x_i|Y=j, \theta) = \prod_i \sum_j w_j f(x_i; \theta_j)$

log-likelihood: $\ell(\theta) = \sum_i \log \sum_j w_j f(x_i; \theta_j)$

**E-Step:** Compute posterior=expected labels=responsibilities: $\gamma_j^{(t)}(x_i) = P(Y=j|x_i, \theta)$

$\gamma_j^{(t)}(x_i) = \frac{w_j^{(t)} f(x_i; \theta_j^{(t)})}{\sum_{l} w_l^{(t)} f(x_i; \theta_l^{(t)})}$, $\frac{P(x_i|Y=j, \theta)P(Y=j|\theta)}{P(x_i;\theta)}$

**Naive Bayes classifier with n binary features:** $\theta_{y,i} = P(X_i=1|Y=y) = \frac{\sum_{j:y_j=y} \mathbb{1}[x_{j,i}]}{\sum_{j:y_j=y}}$

**Labelled and unlabelled data:** $y_{1:n} = $ unlabelled! E-Step: for unlabelled: log-likelihood. $y_j^{(t)}(x_i) = P(X_{1:n}=x_{1:n}|y=j) = $ random or as in $k$-means+, $\Sigma_j = $ data. E-Step with $y_j^{(t)}(x_i)$ and $\mu_j^{(t)}$

**Other method:** Define $Q(y;x) = P(Y_{1:n}|x_{1:n}) = R(\theta)$ but $y_c$ exchanged to $E_Q \Sigma Y$

**M-Step:** $\theta^{(t)} = \arg\max_\theta \sum_m \sum_{j} \gamma_j^{(t)}(x_i) \ln P(x_i, y_i=j; \theta) \Rightarrow$ Expected log-likelihood $E_Q[R(\theta)]$

$w_j^{(t)} = \frac{1}{n} \sum_i \gamma_j^{(t)}(x_i)$, $\mu_j^{(t)} = \frac{\sum_i \gamma_j^{(t)}(x_i) x_i}{\sum_i \gamma_j^{(t)}(x_i)}$, $\Sigma_j^{(t)} = \frac{\sum_i \gamma_j^{(t)}(x_i)(x_i-\mu_j^{(t)})(x_i-\mu_j^{(t)})^T}{\sum_i \gamma_j^{(t)}(x_i)}$

$+ V^2 I$ via CV

**MAP** with (Wishart) prior via CV

Adversaries guaranteed to conv to local opt. Initialization! $y_j^{(t)}(x_i) = P(X=x_i; \theta_j) = \frac{w_j f(x_i; \theta_j) + \lambda \mathbb{1}[y_c = j]}{\sum_i f(x_i; \theta_i)}$ for labeled, $w_j$ normal