

Introduction to Machine Learning Summary

Supervised Learning

Regression

Dataset: $D = \{(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)\} \rightarrow X = [\vec{x}_1, \dots, \vec{x}_n]^T \in \mathbb{R}^{n \times d}$, $\vec{y} = [y_1, \dots, y_n]^T \in \mathbb{R}^n$
Empirical Risk: $R_D(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(\vec{x}_i)) \rightarrow \hat{f}_{\text{EM}} = \arg\min_f R_D(f)$
Population Risk: $R(f) = \mathbb{E}_{(\vec{x}, y) \sim P} [\ell(y, f(\vec{x}))] \rightarrow f^*$ minimizer of $R(f)$

Least Squares Regression:

Objective: $\hat{w} = \arg\min_w \sum_{i=1}^n (y_i - \vec{w}^T \vec{x}_i)^2 = \arg\min_w \|\vec{y} - X\vec{w}\|_2^2$
CF $\rightarrow n > d \rightarrow$ one best solution $\rightarrow \hat{w} = (X^T X)^{-1} X^T \vec{y}$
Samples $n < d \rightarrow$ many exact solutions $\rightarrow \hat{w} = X^T (X X^T)^{-1} \vec{y}$

Ridge Regression:

Uses regularization to encourage small weights
Objective: $\hat{w} = \arg\min_w \|\vec{y} - X\vec{w}\|_2^2 + \lambda \|\vec{w}\|_2^2$ λ_2 Regularization
Closed form solution: $\hat{w} = (X^T X + \lambda I)^{-1} X^T \vec{y}$

Lasso Regression:

Uses regularization to encourage some $w_i = 0$
Objective: $\hat{w} = \arg\min_w \|\vec{y} - X\vec{w}\|_2^2 + \lambda \|\vec{w}\|_1$ λ_1 regularization
No CF in general, intrinsic/embedded method for feature selection

Bias-Variance Tradeoff: $\lambda \rightarrow 0$: Variance \uparrow , Bias \downarrow
 $\lambda \rightarrow \infty$: Variance \downarrow , Bias \uparrow

Model Selection & Validation

True Risk $R(\vec{w})$ would lead to ideal weights $\vec{w}^* = \arg\min_{\vec{w}} R(\vec{w})$, but is unknown \rightarrow estimate over empirical risk $R_D(\vec{w})$ based on D

Empirical Risk Minimization:

$\hat{f}_{\text{EM}} = \arg\min_f R_D(f)$
At \hat{f} the empirical risk $R_D(\hat{f})$ systematically underestimates the population risk: $\mathbb{E}[R_D(\hat{f})] = \mathbb{E}[R(\hat{f})] = R(\hat{f})$
For $n \rightarrow \infty$: $R_D(\hat{f}) \rightarrow R(\hat{f})$ not
 \rightarrow on validation set: $\mathbb{E}_D[R_D(\hat{f})] = R(\hat{f})$

Cross Validation:

(avoids bias with independent validation sets)
K-Fold CV: $D = [D_1 | D_2 | \dots | D_K]$
1. Training set is $k-1$ folds, Validation Set 1 fold
2. Validation set gets varied and then the prediction error is the average of validation errors
 \Rightarrow smaller bias $R(\hat{f}_{\text{EM}}) - R(\hat{f}_{\text{CV}})$
bigger variance in $R(\hat{f}_{\text{EM}}) - R(\hat{f}_{\text{CV}})$
smaller variance in estimating avg. population risk $\mathbb{E}_D[R(\hat{f}_{\text{CV}})]$

Optimization

Gradient Descent: (minimizes objective function $R(\vec{w})$)

1. Initialization of random $\vec{w}_0 \in \mathbb{R}^d$
2. Update $\vec{w}^{(t+1)} = \vec{w}^{(t)} - \eta \nabla_{\vec{w}} L(\vec{w}^{(t)})$
Convergence: $\|\vec{w}_t - \vec{w}^*\|_2 \leq \|\vec{I} - \eta X^T X\|_{\text{op}}^t \cdot \|\vec{w}_0 - \vec{w}^*\|_2$
 \rightarrow Exponential/Linear Convergence if $\rho(\vec{I} - \eta X^T X) = \frac{K-1}{K}$
Iterations for ϵ -close to \vec{w}^* : $T \geq \frac{\log(\epsilon/\epsilon_0)}{\log(K/(K-1))}$
Optimal Stepsize $\eta_{\text{opt}} = \frac{\lambda_{\min}(X^T X)}{\lambda_{\max}(X^T X) + \lambda_{\min}(X^T X)}$, $K = \frac{\lambda_{\max}(X^T X)}{\lambda_{\min}(X^T X)}$
 $\rightarrow T \geq O(K \log(1/\epsilon)) \rightarrow$ large $K \rightarrow$ low convergence
GD is less computationally complex than CF for big problem

Minibatch GD: (\vec{w} with subset of samples S)

1. Initialization of random $\vec{w}_0 \in \mathbb{R}^d$
2. Update $\vec{w}^{(t+1)} = \vec{w}^{(t)} - \eta \frac{1}{|S|} \sum_{(\vec{x}, y) \in S} \nabla_{\vec{w}} L(\vec{x}, \vec{w}^{(t)})$
Stochastic Gradient Descent for $|S|=1$
 \rightarrow only one point $(\vec{x}, y) \in D$ is chosen randomly
 \rightarrow converges if $\sum \eta_t = \infty$ & $\sum \eta_t^2 < \infty$ (e.g. $\eta_t = \frac{1}{t}$)
 \rightarrow Batchsize $\downarrow \rightarrow$ Variance \uparrow , memory \downarrow , (escaping saddlepoint) \uparrow
SGD with momentum: avoids being stuck in local optima & oscillations
 $\rightarrow v_t = \eta \nabla_{\vec{w}} L(\vec{w}^{(t)}) + \gamma v_{t-1}$, $\vec{w}^{(t+1)} = \vec{w}^{(t)} + v_t$

Classification

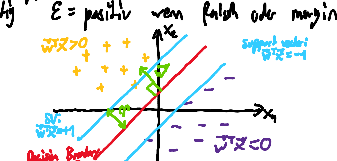
Binary Classification: (minimizing number of mistakes)

Labeling: $Y = \{+1, -1\}$ with $\hat{y} = \text{sign}(f(\vec{x}))$
Objective: $\hat{w} = \arg\min_w \sum_{i=1}^n \ell_{\vec{x}}(\vec{x}_i, \vec{w})$ $R(\hat{f}) = P[\hat{y} \neq \text{sign}(f(\vec{x}))]$
Loss Functions:
• 0-1 $\ell_{0-1} = \mathbb{1}[\hat{y} \neq \text{sign}(f(\vec{x}))]$
• Perceptron $\ell_p = \max\{0, \vec{y} \cdot \vec{w}\}$ (no optimization with explicit)
• Hinge $\ell_h = \max\{0, 1 - \vec{y} \cdot \vec{w}\}$ (close to 0-1)
• Logistic $\ell_{\log} = \log(1 + \exp(-\vec{y} \cdot \vec{w}))$ (least squares)
• Exponential $\ell_{\exp} = \exp(-\vec{y} \cdot \vec{w})$

Logistic Regression: (uses logistic loss)

Objective: $\hat{w} = \arg\min_w \sum_{i=1}^n \log(1 + \exp(-\vec{y}_i \cdot \vec{w}))$
GD: $\vec{w}_{t+1} = \vec{w}_t + \eta \frac{1}{n} \sum_{i=1}^n \vec{x}_i \vec{y}_i / (1 + \exp(\vec{y}_i \cdot \vec{w}_t))$

Support Vector Machine: (uses hinge loss to encourage margin)

Objective: $\hat{w} = \arg\min_w \sum_{i=1}^n \max\{0, 1 - \vec{y}_i \cdot \vec{w}\} + \lambda \|\vec{w}\|_2^2$
Margin(\vec{w}) = $\min_i \vec{y}_i \cdot \vec{w}$ \rightarrow min distance = $\frac{\text{margin}(\vec{w})}{\|\vec{w}\|_2}$
• Hard Margin SVM: $\min \|\vec{w}\|_2$ s.t. $\vec{y}_i \cdot \vec{w} \geq 1 \forall i \in \{1, 2, \dots, n\}$
 \rightarrow for linear separable Data
• Soft Margin SVM: $\min \frac{1}{2} \|\vec{w}\|_2^2 + C \cdot \sum_{i=1}^n \max\{0, 1 - \vec{y}_i \cdot \vec{w}\}$
 \rightarrow for not linear separable Data
 $\rightarrow \epsilon = 0$ max margin ϵ = positive von false oder margin


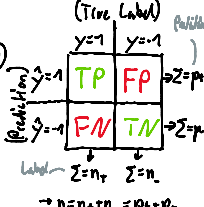
Multiclass Classification:

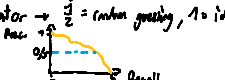
Labels: $Y = \{1, \dots, K\}$
On a All: 1. Train K binary classifiers separately
 \rightarrow Current Labels $\rightarrow m$ Rel $\rightarrow \vec{1}$
2. Apply softmax $\Rightarrow \sigma(\vec{f}(\vec{x})) = \frac{e^{\vec{f}(\vec{x})}}{\sum_{j=1}^K e^{\vec{f}_j(\vec{x})}}$
3. For test point: $\hat{y} = \arg\max_j \sigma(\vec{f}_j(\vec{x})) = \arg\max_j \vec{f}_j(\vec{x})$
Decision boundary for f_i & f_j is $f_i = f_j$ their slope

Hypothesis Testing/Other Metrics:

Pit. class as Null hypothesis ($y=0$) if critical to profit correctly when true
• reject $H_0(x) \Leftrightarrow \hat{y}=1$ if $f(x) > J$ used to control significance level, FPR
• accept $H_0(x) \Leftrightarrow \hat{y}=0$ if $f(x) < J$

Accuracy = $\frac{TP+TN}{n} = P(\hat{y}=y)$
Recall/TPR = $\frac{TP}{TP+FN} = \frac{TP}{n_+} = P(\hat{y}=1|y=1)$
FDR = $1 - \text{Precision} = \frac{FP}{FP+TP} = P(\hat{y}=1|\hat{y}=1)$
Precision = $\frac{TP}{TP+FP} = P(y=1|\hat{y}=1)$
FPR = $\frac{FP}{FP+FN} = \frac{FP}{n_-} = P(\hat{y}=1|y=-1)$
F1 score = $\frac{2 \cdot \text{Precision} \cdot \text{Recall}}{2 \cdot \text{Precision} + \text{Recall}}$



ROC curve: varying γ for fixed $\hat{f} \rightarrow$ plot TPR vs FPR
 $\rightarrow \gamma$ small \rightarrow TPR \uparrow , FPR \uparrow vs γ large \rightarrow TPR \downarrow , FPR \downarrow
 \rightarrow Area under curve: indicator $\rightarrow \frac{1}{2}$ = random guessing, 1 = ideal
Precision vs. Recall curve: 

Non Parametric Methods

k-Nearest Neighbors: (no training, inference during test time $\rightarrow O(nd)$)

1. Given dataset D , pick k and distance metric d
2. For a given \vec{x} , find k closest $\vec{x}_1, \dots, \vec{x}_k$ in D
3. Output majority vote/average of y_1, \dots, y_k
 \rightarrow sensitive to k , pick k using CV, erratic in high dim., logn needed

Decision Trees:

1. Training set $D \rightarrow$ each node v contains $D(v) \in D$
2. Parent node split $D(v)$ with chosen rule ($>$, $<$)
3. repeat until uncertainty $H(S)$ is low

Kernels

Feature expansion needs to be avoided \rightarrow optimal hyperplane lies in span of data $\rightarrow \vec{w}$ has fixed structure
Feature expansion: $\vec{x} \in \mathbb{R}^d$, $\phi(\vec{x}) \in \mathbb{R}^p \rightarrow p = \binom{d+m}{m} \approx O(d^m)$
 $d=2, m=1 \rightarrow (1, x_1, x_2, x_1^2, x_2^2)$
 $d=2, m=2 \rightarrow (1, x_1, x_2, x_1^2, x_2^2, x_1 x_2)$

Kernel Trick:

$\vec{\Phi} = \sum_{i=1}^n \vec{y}_i \vec{x}_i^T = \sum_{i=1}^n \alpha_i \phi(\vec{x}_i) = \vec{\Phi}^T \vec{x}$ with $\vec{\Phi} = \begin{bmatrix} -\phi(\vec{x}_1) \\ \vdots \\ -\phi(\vec{x}_n) \end{bmatrix}$
Kernelization: $\vec{x} \rightarrow \vec{\Phi}$, $\vec{x} \rightarrow \vec{\Phi}(\vec{x})$
Kernel Function: $k(\vec{x}, \vec{x}') = \vec{\Phi}(\vec{x})^T \vec{\Phi}(\vec{x}') = \langle \vec{\Phi}(\vec{x}), \vec{\Phi}(\vec{x}') \rangle$
Kernel/gram Matrix: $K = \vec{\Phi} \vec{\Phi}^T = \begin{bmatrix} k(\vec{x}_1, \vec{x}_1) & \dots & k(\vec{x}_1, \vec{x}_n) \\ \vdots & \ddots & \vdots \\ k(\vec{x}_n, \vec{x}_1) & \dots & k(\vec{x}_n, \vec{x}_n) \end{bmatrix}$

Properties:

- Symmetry: $k(\vec{x}, \vec{x}') = k(\vec{x}', \vec{x})$
- K is positive semidefinite \rightarrow Eigenvalues ≥ 0

Composition of Kernel Functions:

- $k(\vec{x}, \vec{x}') = k_1(\vec{x}, \vec{x}') + k_2(\vec{x}, \vec{x}')$ • $k(\vec{x}, \vec{x}') = k_1(\vec{x}, \vec{x}') \cdot k_2(\vec{x}, \vec{x}')$
- $k(\vec{x}, \vec{x}') = c \cdot k_1(\vec{x}, \vec{x}')$ for $c > 0$ • $k(\frac{\vec{x}}{c}, \frac{\vec{x}'}{c}) = k_1(\vec{x}, \vec{x}')$
- $k(\vec{x}, \vec{x}') = f(k_1(\vec{x}, \vec{x}'))$ if polynomial with pos. coeff. or exp
- $k(\vec{z}, \vec{z}') = k_1(\vec{v}(\vec{z}), \vec{v}(\vec{z}'))$ $\vec{v}: \mathbb{Z} \rightarrow \mathbb{X}$ • $k(\frac{\vec{x}}{c}, \frac{\vec{x}'}{c}) = k_1(\vec{x}, \vec{x}') k_2(\vec{y}, \vec{y}')$

Important Kernel Functions:

- Constant: $k(\vec{x}, \vec{x}') = c \geq 0$ • Laplacian $k(\vec{x}, \vec{x}') = \exp(-\frac{\|\vec{x} - \vec{x}'\|_1}{h})$
- Linear: $k(\vec{x}, \vec{x}') = \vec{x}^T \vec{x}'$
- Gaussian/RBF: $k(\vec{x}, \vec{x}') = \exp(-\frac{\|\vec{x} - \vec{x}'\|_2^2}{2\sigma^2})$ \leftarrow infinite dimensional
- Polynomial (Deg. m): $k(\vec{x}, \vec{x}') = (\vec{x}^T \vec{x}')^m$ \leftarrow (deg. m) = $O(d^m)$ for Deg. m monomials
- Polynomial (Up to Deg. m): $k(\vec{x}, \vec{x}') = (1 + \vec{x}^T \vec{x}')^m$

Neural Networks

Parameterize feature maps ϕ & optimize over parameters
 \rightarrow results in a non-convex function for multilayer networks

Activation Functions: (\rightarrow map the data (non-linearly))

- Linear: $\varphi(z) = z$, $\varphi(z) = 1$ (linear)
- Sigmoid: $\varphi(z) = \frac{1}{1 + \exp(-z)}$, $\varphi(z) = \varphi(z) \cdot (1 - \varphi(z))$ (soft edges)
- Tanh: $\varphi(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$, $\varphi(z) = 1 - \varphi(z)^2$
- ReLU $\varphi(z) = \max\{0, z\}$, $\varphi(z) = \begin{cases} z & z \geq 0 \\ 0 & z < 0 \end{cases}$ (rare edges)

Forward Propagation: (predict)

- Input Layer O : $\vec{z}^{(0)} = [\vec{x}, 1]$ initialize with x values or 1
 - Hidden Layer L : $\vec{z}^{(L)} = W^{(L)} \vec{z}^{(L-1)}$ sum weighted outputs of $L-1$
 $\vec{a}^{(L)} = [\varphi(\vec{z}^{(L)})]$ output sum with act. function
 - Output Layer L : $\vec{f} = W^{(L)} \vec{z}^{(L-1)}$ sum weighted outputs of $L-1$
- \rightarrow Predict $\vec{y}_i = f_i$ for regression & $\vec{y} = \arg\max_j f_j$ for classification

Backpropagation: (adjust weights)

- Output Layer: Error: $\delta^{(L)} = \nabla_{\vec{f}} L(\vec{f}) = [\frac{\partial L}{\partial f_1}, \dots, \frac{\partial L}{\partial f_p}]$
• Gradient: $\nabla_{\vec{z}^{(L)}} L(\vec{y}, \vec{z}; W) = \delta^{(L)} \cdot \vec{1}$
 - Hidden Layer: Error: $\delta^{(L)} = \frac{\partial L}{\partial \vec{z}^{(L)}} = \varphi'(\vec{z}^{(L)}) \odot (W^{(L+1)T} \delta^{(L+1)})$ (minimize multiplication)
• Gradient: $\nabla_{\vec{z}^{(L)}} L(\vec{y}, \vec{z}; W) = \delta^{(L)} \cdot \vec{1}$
- \rightarrow Take derivative of loss function up to weight $\frac{\partial}{\partial x} g(\vec{x}) \frac{\partial \vec{x}}{\partial x} = \frac{\partial f}{\partial x}$

Exploding/Disappearing Gradient, Weight initialization:

Exploding gradient: $\|\nabla_{\vec{w}} L\| \rightarrow \infty$, Vanishing gradient: $\|\nabla_{\vec{w}} L\| \rightarrow 0$
By using ReLU and keeping the variance of the weights a constant across layers, vanishing/exploding gradients can be avoided
For a good weight initialization we need to set $\sigma^2 = \frac{2}{n}$
 \rightarrow Glorot (fan) $\rightarrow w_{ij} \sim N(0, \frac{1}{n_{in}})$, He (ReLU) $\rightarrow w_{ij} \sim N(0, \frac{2}{n_{in}})$
 n_{in} = number of units in previous layer, n_{out} = number of units in next layer

Regularization: (avoid overfitting)

- Add penalty term: $+ \lambda ||W||_2^2$
- Early Stopping: Stop training with GD before convergence
- Dropout: randomly dropout units for each Health with prob. = p
- Batch Normalization: Normalization of inputs to $N(0,1)$ with $\mu = \frac{1}{N} \sum_{i=1}^N x_i$, $\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2 \rightarrow \hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta$ Shift

Convolutional Neural Networks: (m-f) Parameters

- Stride s: Number of Pixels Filter moves
- Padding p: Number of "0" Pixels added to side of image
- Filter & Kernel: Filter is a concatenation of Kernels (dim(F) = dim(K) * s)
- Output Dimension: Applying m different fxf filters to an m x m image results in an Output dim = L x L x m with $L = \frac{n - f + 1}{s} + 1$
- Pooling: Reduces promotes by only considering average/maximum of multiple units

Unsupervised Learning

Clustering

k - Means Clustering: (Find centres that minimize avg. squared distances)

- Objective: $\hat{\mu} = \arg \min_{\mu} R(\hat{\mu}) = \arg \min_{\mu} \sum_{i=1}^n \min_{j \in \{1, \dots, k\}} ||x_i - \mu_j||_2^2$ (non-convex)
- ↳ Lloyd's Heuristic:
1. Initialize cluster centres $\mu^{(0)} = [\mu_1^{(0)}, \dots, \mu_k^{(0)}]$
 2. While not converged: (no more changes)
 - i) Assign points to closest centre: $z_i = \arg \min_{j \in \{1, \dots, k\}} ||x_i - \mu_j^{(t-1)}||_2$ point i to number of centre i
 - ii) Update centres: $\mu_j^{(t)} = \frac{1}{n_j} \sum_{i: z_i=j} x_i$
- converges to local optimum with cost per iteration = $O(nkd)$

k - Means ++: (Initialization for k-Means)

1. Start with random data point as first cluster centre
 2. Add 2-k centres randomly proportional to squared distance to closest centre = $\frac{2}{n} \sum_{i=1}^n ||x_i - \mu^{(0)}||_2^2$
- To find the optimal number of clusters find the "elbow" in \hat{R} vs. k plot or add regularizer (L1)
- Expected Loss of k-Means ++: $O(\log(n)) \cdot \min_{\mu} R(\mu)$

Dimension Reduction - Principal Component Analysis (PCA)

PCA: (two dimensional representation in direction with highest variance)

- Given: • Centered Data with $\mu = \sum_{i=1}^n x_i = 0$
- Empirical Covariance Matrix $\Sigma = \frac{1}{n} \sum_{i=1}^n x_i x_i^T = \frac{1}{n} \sum_{i=1}^n \begin{bmatrix} x_{i1}^2 & x_{i1}x_{i2} \\ x_{i1}x_{i2} & x_{i2}^2 \end{bmatrix}$
- Objective: $\hat{W} = \arg \min_{W} \sum_{i=1}^n ||Wz_i - x_i||_2^2$
- Solution: In principle eigenvectors of $\Sigma \rightarrow \hat{W} = [v_1, \dots, v_k] \in \mathbb{R}^{d \times k}$ with $z_i = W^T x_i$, $\hat{Z} = \sum_{i=1}^n \lambda_i v_i v_i^T$
- minimizes Euclidean rec. error, perfect reconstruction for $k > n$ & $k = d$

Kernel PCA: (can be used to discover non-linear feature maps h(x))

- Objective: $\hat{\alpha} = \arg \min_{\alpha} ||\sum_{i=1}^n \alpha_i^T K(i, j) - y||_2^2$
- Solution: In principle eigenvectors of $K \rightarrow \hat{\alpha} = [\alpha_1^{(1)}, \dots, \alpha_n^{(1)}]^T$ with $\alpha_i^{(1)} = \frac{1}{\lambda_i} \langle v_i, \tilde{v} \rangle$ and $K = \sum_{i=1}^n x_i x_i^T = \Phi \Phi^T$
- Projection: $z_i = W^T x_i = \sum_{j=1}^n \alpha_j^T k(x_i, x_j)$

Auto encoder

Objective: Learn $f(x; W) = f_{dec}(f_{enc}(x; W_{enc}); W_{dec}) \approx x$ with $\hat{W} = \arg \min_{W} \sum_{i=1}^n ||x_i - f(x_i; W)||_2^2$

→ Neural Network with $k < d$ hidden layer units

→ equivalent to PCA if activation function is identity

Probabilistic Modeling

Probability Formulas

Exp. Value: $E[X] = \sum_{i=1}^n x_i \cdot P(x_i)$, Variance: $Var(X) = E[X^2] - E[X]^2$

Bayes Rule: $P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$, $P(A|B) = \frac{P(B|A) \cdot P(A)}{\sum_{i=1}^n P(B|A_i) \cdot P(A_i)} = \frac{P(A, B)}{P(B)}$

MLE: $\hat{\theta}_{MLE} = \arg \max_{\theta} P(\theta; D) = \arg \max_{\theta} \prod_{i=1}^n P(x_i; \theta) = \arg \max_{\theta} - \sum_{i=1}^n \log(P(x_i; \theta))$

MAP: $\hat{\theta}_{MAP} = \arg \max_{\theta} P(\theta; D) = \arg \max_{\theta} \frac{P(\theta) \cdot P(D|\theta)}{P(D)} = \arg \max_{\theta} \frac{P(\theta)}{P(D)} \cdot \log(P(D|\theta))$

- Normal Distribution $N(x; \mu, \Sigma = \sigma^2) \rightarrow P(x=x_i) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x_i - \mu)^2}$
- Bernoulli Distribution $Be(p) \rightarrow P(x=x_i) = p^{x_i} (1-p)^{1-x_i}$, $x_i \in \{0, 1\}$
- Binomial Distribution $Bi(n, p) \rightarrow P(x=x_i) = \binom{n}{x_i} p^{x_i} (1-p)^{n-x_i}$, $x_i \in \{0, 1, \dots, n\}$
- Poisson Distribution $Po(\lambda) \rightarrow P(x=x_i) = \frac{\lambda^{x_i}}{x_i!} e^{-\lambda}$, $x_i \in \{0, 1, \dots, \infty\}$

Decision Theory

Decision Rule:

A decision rule $f: X \rightarrow A$ takes an observable input $x \in X$ and outputs an action $f(x)$ in some action space A which causes outcomes $(f(x), y)$

→ optimal decision rule f^* minimizes $E[L(f(x), y)]$ and hence $R(f)$

↳ optimal action: $a^*(x) = \arg \min_{a \in A} E_{y \sim p(y|x)} [L(a, y)]$ need $P(y|x)$

Asymmetric Losses: $L(a, y) = J \cdot \max(y - a, 0) + (1 - J) \cdot \max(a - y, 0) = (y - a)(J - J_0 + J_0 \cos \theta)$

0-1 Loss: $L(a, y) = C_0 \mathbb{1}_{a=y} + C_1 \mathbb{1}_{a \neq y} = C_0 \mathbb{1}_{a=y} + C_1 (1 - \mathbb{1}_{a=y}) = C_1 + (C_0 - C_1) \mathbb{1}_{a=y}$

Discriminative & generative Models

Discriminative Models: Y is an uncertain consequence of $X \rightarrow P(y|x)$

- Food ingredients (X) → Healthy? (Y)
- Model: $P(x, y; \theta) = P(y|x; \theta) \cdot P(x)$
- Regression or Classification

Generative Models: X is an uncertain consequence of $Y \rightarrow P(x, y)$

- Disease (Y) → Symptoms (X)
- Model: $P(x, y; \theta) = P(x|y; \pi) \cdot P(y; \pi) = \theta(y, \pi)$
- More common for classification (harder problem)

Gaussian Bayes Classifier (GBC)

Class Prior: $\forall y \in \{1, 2\}, P(y=y; \pi) = \pi_y = \pi_j$

Feature Distribution: $X|Y \sim N(x; \mu_y, \Sigma_y) = \frac{1}{(2\pi)^{d/2} |\Sigma_y|^{1/2}} e^{-\frac{1}{2} (x - \mu_y)^T \Sigma_y^{-1} (x - \mu_y)}$

MLE's: $\hat{\pi}_j = \frac{\text{Count}(Y=j)}{n}$

- $\hat{\mu}_j = \frac{\sum_{i: y_i=j} x_i}{\text{Count}(Y=j)}$
- $\hat{\Sigma}_j = \frac{\sum_{i: y_i=j} (x_i - \hat{\mu}_j)(x_i - \hat{\mu}_j)^T}{\text{Count}(Y=j) - 1}$

Gaussian Naive Bayes: features independent $\rightarrow \Sigma_y = \text{diag}(\sigma_1^2, \dots, \sigma_d^2)$

Quadratic Discriminant Analysis: π_j is uniform

→ Linear Discriminant Analysis: $\Sigma_y = \Sigma$ for all y

→ Fisher's LDA: $c=2, \mu_1 = \mu_2 = \mu, \Sigma_1 = \Sigma_2$

Frequentists Inference

Bootstrap: No assumptions on model needed → estimate dist. via sample

Estimate $R(f_0) \approx E_0[R(\hat{f}_0)]$ with $\hat{f}_0(x_i) = \frac{1}{n} \sum_{j=1}^n I(x_j = x_i) y_j$

Principle: 1. Sample $D = \{x_i, y_i\}_{i=1}^n$ i.i.d. from $P(X \times Y)$ uniformly

2. Bootstrap: For $b=1, \dots, B$ D_b consists of uniform samples $x_i^{(b)}$ with replacement from D

3. Compute $\hat{\theta}_b = \frac{1}{n} \sum_{i=1}^n x_i^{(b)}$ for every $D_b \rightarrow$ Distribution for $\hat{\theta}$

Justification: Instead of sampling D_b with replacement, use K -fold CV

Gaussian Mixture Models (GMM)

Combination of Gaussians: $P(x|\theta) = \sum_{i=1}^K \pi_i N(x; \mu_i, \Sigma_i)$ with $\pi_i \geq 0, \sum_{i=1}^K \pi_i = 1$

↳ $\log(P(x|\theta)) = \log(\sum_{i=1}^K \pi_i N(x; \mu_i, \Sigma_i)) = \log(\sum_{i=1}^K \pi_i \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} e^{-\frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)})$

Hard EM: assign fixed labels

E-Step: Predict most likely class for each data point

→ $z_i^{(t)} = \arg \max_{j \in \{1, \dots, K\}} P(z_i=j; \theta^{(t-1)}) = \arg \max_{j \in \{1, \dots, K\}} \pi_j N(x_i; \mu_j, \Sigma_j)$

M-Step: MLE for now complete Data $D^{(t)} = \{(x_i, z_i^{(t)})\}_{i=1}^n$

→ $\theta^{(t)} = \arg \max_{\theta} P(D^{(t)}|\theta)$ fit $\hat{\mu}_j$ & $\hat{\Sigma}_j$ as with GBC

→ problematic if 2 clusters of same class overlapping

Soft EM: assign probabilities

E-Step: Calculate Responsibilities

$Q(\theta, \theta^{(t-1)}) = E_{z \sim \pi} [\log(P(x_{i1}, \dots, x_{in}|\theta))] x_{in}, \theta^{(t-1)}$

→ $\pi_j^{(t)}(x_i) = P(z_i=j; \theta^{(t-1)}) = \frac{\pi_j N(x_i; \mu_j, \Sigma_j)}{\sum_{k=1}^K \pi_k N(x_i; \mu_k, \Sigma_k)}$

M-Step: Fit clusters to weighted data point by maximizing log-likelihood

$\theta^{(t)} = \arg \max_{\theta} Q(\theta, \theta^{(t-1)}) = \arg \max_{\theta} \sum_{i=1}^n \sum_{j=1}^K \pi_j^{(t)}(x_i) \log(\pi_j N(x_i; \mu_j, \Sigma_j))$

→ $\mu_j^{(t)} = \frac{1}{n_j} \sum_{i=1}^n \pi_j^{(t)}(x_i) x_i$, $\pi_j^{(t)} = \frac{\sum_{i=1}^n \pi_j^{(t)}(x_i)}{\sum_{i=1}^n \sum_{j=1}^K \pi_j^{(t)}(x_i)}$

$\Sigma_j^{(t)} = \frac{\sum_{i=1}^n \pi_j^{(t)}(x_i) (x_i - \mu_j^{(t)})(x_i - \mu_j^{(t)})^T}{\sum_{i=1}^n \pi_j^{(t)}(x_i)} + \gamma^{-1} I$

without prior

Initialization:

Weights	Means	Variances
uniform distribution	randomly/cluster	leptokurtic with empirical variance

→ either k or γ^2 is selected with CV module log-likelihood on testset

Combined GMM's:

	Spherical	Diagonal	Tied	Full
Parameters	k	$k \cdot d$	$\frac{d(d+1)}{2}$	$\frac{d(d+1)}{2} + k$

Degeneracy: GMM chooses $k=n$ (1 Gaussian band each datapoint), thus $\Sigma \rightarrow 0$ which leads to overfitting. Can be avoided by without prior (π_i)

Anomaly Detection: $P(x|\theta) < J \rightarrow$ detect x_i as anomaly (outlier FPR via J)

Generative Adversarial Networks (GAN)

2 Neural networks are trained. The generator tries to produce realistic data out of noise and the discriminator tries to spot fake data. Training leads to a saddle point between both.

Objective: $\min_G \max_D E_{x \sim p_G} [\log(D(x; w_D))] + E_{z \sim p_Z} [\log(1 - D(z; w_D))] = \min_G \max_D [E_{x \sim p_G} \log(D(x; w_D))] + E_{z \sim p_Z} [\log(1 - D(z; w_D))]$

Performance Metric: Daily log- $\log \rightarrow D_h(w_D, w_G) = \frac{1}{n} \sum_{i=1}^n \log(D(x_i; w_D)) - \frac{1}{n} \sum_{i=1}^n \log(1 - D(z_i; w_D))$

Optimal Discriminator: $D_h^*(x) = \frac{P_G(x) + P_Z(x)}{P_G(x) + P_Z(x)}$ → splits every \rightarrow fake data with $P = \frac{P_G(x)}{P_G(x) + P_Z(x)}$

Optimal generator: $P_G(x) = P_{min}(x) \rightarrow D_h^*(x) = \frac{1}{2}$

Other Formulas

Linear Algebra

Matrix composition: $XX^T = [x_1, x_2, x_3] \begin{bmatrix} x_1^T \\ x_2^T \\ x_3^T \end{bmatrix} = \begin{bmatrix} x_1^T x_1 & x_1^T x_2 & x_1^T x_3 \\ x_2^T x_1 & x_2^T x_2 & x_2^T x_3 \\ x_3^T x_1 & x_3^T x_2 & x_3^T x_3 \end{bmatrix}$

Transpose: $S = A^T \rightarrow S^T = A$

Scalar Product: $\langle \vec{u}, \vec{v} \rangle = \vec{u}^T \vec{v} = \sum_{i=1}^n u_i v_i$, $\langle \vec{v}, \vec{v} \rangle = ||\vec{v}||_2^2 = \vec{v}^T \vec{v}$

Trace: $\text{tr}(A) = \sum_{i=1}^n a_{ii}$, $\text{tr}(\vec{v} \vec{v}^T) = \text{tr}(\vec{v}^T \vec{v})$

Orthogonal Matrices: $A^T = A^{-1} \rightarrow A^T A = I$ (columns vectors $||v_i||_2=1$ & orthogonal $\rightarrow v_i^T v_j = 0$)

Inverse: $A^{-1} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{\det(A)} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$ only when squared and $\det(A) \neq 0$

Rank Nullity Theorem: $\text{rank}(A) + \text{null}(A) = \dim(\text{vec}(A)) = \text{rank}(A) + n$

Eigenvalues and vectors: $A \vec{v} = \lambda \vec{v}$, $\vec{v} \in EV$, $\lambda \in EW$

Linear Differential Calculus: $\frac{dA}{dx} = A$, $\frac{dA^T}{dx} = A^T$

$\frac{d(A^T A)}{dx} = (A^T + A) \frac{dA}{dx}$

SVD: $A = U \Sigma V^T$ for $m \geq n$ (else $A^T = U^T V^T$)

- $V = [\vec{v}_1, \dots, \vec{v}_n]$ \vec{v}_i is EV of $A^T A$
- $S = \text{diag}(s_1, \dots, s_n) + 0$ $s_i = \sqrt{\lambda_i}$ λ_i EW of $A^T A$
- $U = [\vec{u}_1, \dots, \vec{u}_m]$ $\vec{u}_i = \frac{1}{\lambda_i} A \vec{v}_i$ (orthogonal)

L2 - Norm: $||x||_2 = \sqrt{x_1^2 + \dots + x_n^2}$

Complexity: $A \in \mathbb{R}^{m \times n}, B \in \mathbb{R}^{n \times m}, C \in \mathbb{R}^{m \times m} \rightarrow AB = O(nmp)$ $C^{-1} = O(n^3)$

Analysis

Differentiation Rules: $f(x)g(x) = f'(x)g(x) + f(x)g'(x)$

$f(g(x)) = f'(g(x))g'(x)$ ← right side important

Convexity: $f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y)$, $0 \leq \theta \leq 1$

→ f convex if $f''(x) \geq 0 \forall x$

→ f strictly convex if $f''(x) > 0 \forall x$

Local minima of convex & strictly convex functions are global minima

Gradient: $\nabla f(x_0) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}_{x=x_0}$

Minimum if $\nabla f=0$ & HJ pos/pos

Maximum if $\nabla f=0$ & HJ neg/neg

Subgradient: Grad/ent at non-differentiable point of $f(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$

Hessian: $H = \nabla^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}_{x=x_0}$

Point of Observation

→ H positive semi-definite ($\lambda_i \geq 0$) $\Rightarrow f$ convex ($\vec{v}^T H \vec{v} \geq 0 \forall \vec{v}$)

→ H positive definite ($\lambda_i > 0$) $\Rightarrow f$ strictly convex ($\vec{v}^T H \vec{v} > 0 \forall \vec{v}$)