

1 Bayesian Decision Theory

Bernoulli : $P(X) = p_0^x(1 - p_0)^{1-x}$, p_0 is the param.

Estimation of p_0 form $\mathcal{X} = \{x^{(i)}\}_{i=1}^N$:
 $\hat{p}_0 = \frac{\text{\#heads}}{\text{\#tosses}} = \frac{\sum_{i=1}^N x^{(i)}}{N}$

Predict outcome to head if $P_0 > 1/2$, tail otherwise.

Bayes's Rule :

Posterior $P(C|\mathbf{x}) = \frac{\text{likelihood} \cdot \text{prior}}{\text{evidence}} = \frac{p(\mathbf{x}|C_i)P(C_i)}{p(\mathbf{x})}$

Baye for $K \geq 2$ classes:

$P(C_i|\mathbf{x}) = \frac{p(\mathbf{x}|C_i)P(C_i)}{\sum_{k=1}^K p(\mathbf{x}|C_k)P(C_k)}$

Optimal decision: Choose C_i if $P(C_i|\mathbf{x}) = \max_k (C_k|\mathbf{x})$

Losses and Risks :

$R(\alpha_i|\mathbf{x}) = \sum_{k=1}^K \lambda_{ik} P(C_k|\mathbf{x})$,
 α_i the action assigned to class C_i , λ_{ik} loss for α_i if C_k

Optimal α_i if $R(\alpha_i|\mathbf{x}) = \min_k R(\alpha_k|\mathbf{x})$

0-1 loss : $R(\alpha_i|\mathbf{x}) = \sum_{k=1}^K \lambda_{ik} P(C_k|\mathbf{x}) = 1 - P(C_i|\mathbf{x})$,
 $\lambda_{ik} = 1$ if $i \neq k$, 0 otherwise

Reject Option : Loss:

$$\lambda_k = \begin{cases} 1 & \text{if } i = k \\ \lambda & \text{if } i = K + 1 \\ 0 & \text{otherwise} \end{cases}$$

λ is the loss for choosing reject.

Expected risk:

$R(\alpha_i|\mathbf{x}) = \sum_{k=1}^K \lambda P(C_k|\mathbf{x}) = \lambda$ if $i = K + 1$
 $P(C_i|\mathbf{x}) = \frac{p(\mathbf{x}|C_i)P(C_i)}{\sum_{k=1}^K p(\mathbf{x}|C_k)P(C_k)}$ if $i = 1, \dots, K$

Optimal Decision: choose C_i if

$R(\alpha_i|\mathbf{x}) = \min_{1 \leq k \leq K} R(\alpha_k|\mathbf{x}) < R(\alpha_{K+1}|\mathbf{x})$,

Reject otherwise.

Discriminant Functions : choose C_i if $g_i(\mathbf{x}) = \max_k g_k(\mathbf{x})$

$g_i(\mathbf{x}) = -R(\alpha_i|\mathbf{x}) = P(C_i|\mathbf{x}) = p(\mathbf{x}|C_i)P(C_i)$

Two class may define single discriminant functions: $g(\mathbf{x}) =$

$g_1(\mathbf{x}) - g_2(\mathbf{x})$, choose C_i if it greater than zero.

Decision Regions $\mathcal{R}_i = \{\mathbf{x}|g_i(\mathbf{x}) = \max_k g_k(\mathbf{x})\}$

Bayesian Network joint:

$P(X_1, \dots, X_d) = \prod_{i=1}^d P(X_i | \text{parents}(X_i))$

2 Parameter Estimation

Setting Assume data follow a distribution model $\mathcal{X} = \{\mathbf{x}^{(i)}\}$,

$\mathbf{x}^* \sim p(\mathbf{x})$. Assume some parametric form for $p(\mathbf{x}|\theta)$, θ is estimated using \mathcal{X}

Param Approach to classification In Baye's rule for classification, $p(\mathbf{x}|C_i)$ (likelihood) and $P(C_i)$ (prior) need to be estimated from the sample \mathcal{X}

MLE seeks to find θ that makes sampling from $p(\mathbf{x}|\theta)$ as likely as possible

likelihood:

$L(\theta|\mathcal{X}) \equiv p(\mathcal{X}|\theta) = \prod_{i=1}^N p(\mathbf{x}^{(i)}|\theta)$

Principle Component Analysis Projection of \mathbf{x} on the direction of \mathbf{w} : $\mathbf{z} = \mathbf{w}^T \mathbf{x}$

Finding the first principle component \mathbf{w}_1 such that $\text{Var}(\mathbf{z}_1)$ is maximized:

$\text{Var}(\mathbf{z}_1) = \text{Var}(\mathbf{w}^T \mathbf{x}) = E[(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \mu)^2]$

$= \mathbf{w}^T E[(\mathbf{x} - \mu)(\mathbf{x} - \mu)^T] \mathbf{w} = \mathbf{w}^T \Sigma \mathbf{w}$

$\text{Cov}(\mathbf{x}) = E[(\mathbf{x} - \mu)(\mathbf{x} - \mu)^T] = \Sigma$

Lagrangian: $\mathbf{w}_1^T \Sigma \mathbf{w}_1 - \alpha (\mathbf{w}_1^T \mathbf{w}_1 - 1)$

Taking derivative: $\Sigma \mathbf{w}_1 = \alpha \mathbf{w}_1$ (eigenvalue equation)

$\mathbf{w}_1^T \Sigma \mathbf{w}_1 = \alpha \mathbf{w}_1^T \mathbf{w}_1 = \alpha$ (eigenvalue)

We choose the eigenvector with the largest eigenvalue for the variance to be maximum.

Second PC: $\mathbf{w}_2^T \Sigma \mathbf{w}_2 - \alpha (\mathbf{w}_2^T \mathbf{w}_2 - 1) - \beta (\mathbf{w}_2^T \mathbf{w}_1 - 0)$

Derivative: $2\Sigma \mathbf{w}_2 - 2\alpha \mathbf{w}_2 - \beta \mathbf{w}_1 = 0$ (times \mathbf{w}_1^T on both side, $\mathbf{w}_1^T \mathbf{w}_2 = 0$)

Then $\Sigma \mathbf{w}_2 = \alpha \mathbf{w}_2$, \mathbf{w}_2 is the second largest eigenvalue.

Proportion of variance (PoV) Explained: $\frac{\lambda_1 + \dots + \lambda_k}{\lambda_1 + \dots + \lambda_d}$

Factor Analysis Assume latent factors \mathbf{z}_j Sample $\mathcal{X} = \{\mathbf{x}^{(i)}\}$, $E(\mathbf{x}) = \mu$, $\text{Cov}(\mathbf{x}) = \Sigma$

Factors \mathbf{z}_j , $E[\mathbf{z}_j] = 0$, $\text{Var}(\mathbf{z}_j) = 1$

Noise: ϵ_i : $E[\epsilon_i] = 0$, $\text{Var}(\epsilon_i) = \Psi_i$, $\text{Cov}(\epsilon_i, \epsilon_j) = 0$

$\mathbf{x}_i - \mu = \sum_{j=1}^K v_{ij} \mathbf{z}_j + \epsilon_i$, assume $\mu = 0$, v_{ij} are called factor loading

$\text{Var}(\mathbf{x}_i) = \sum_{j=1}^K v_{ij}^2 \text{Var}(\mathbf{z}_j) + \text{Var}(\epsilon_i) = \sum_{j=1}^K v_{ij}^2 + \Psi_i$

Covariance Matrix: $\Sigma = \text{Cov}(\mathbf{V}\mathbf{z} + \epsilon) = \mathbf{V}\mathbf{V}^T + \Psi$

Factor loading: $\text{Cov}(\mathbf{x}, \mathbf{z}) = \mathbf{V}$

Dim reduce: Given \mathbf{S} as the estimator of Σ , we want to find \mathbf{V} and Ψ s.t. $\mathbf{S} = \mathbf{V}\mathbf{V}^T + \Psi$, $\Psi = \text{diag}(\Psi_i)$

Multidimensional Scaling lower dimension preserve pairwise distances.

Sample $\mathcal{X} = \{\mathbf{x}^{(i)} \in \mathbb{R}^d\}_{i=1}^N$

Squared Euclidean distance between point i and j :

$d_{ij}^2 = \sum_{j=1}^d (x_j^{(i)} - x_j^{(j)})^2 = b_{ij} + b_{jj} - 2b_{ij}$

$b_{ij} = \sum_{k=1}^d x_k^{(i)} x_k^{(j)}$ or matrix form $\mathbf{B} = \mathbf{X}\mathbf{X}^T$

Constraint: $\sum_j x_j^{(i)} = 0, \forall i$, define: $T = \sum_{i=1}^N b_{ii}$

Then $\sum_i d_{ij}^2 = T + Nb_{jj}$, $\sum_i \sum_j d_{ij}^2 = d_{jj}^2$

defining: $d_{jj}^2 = \sum_{i=1}^N v_{ij}^2 \text{Var}(\mathbf{z}_j) + \text{Var}(\epsilon_j) = \sum_{i=1}^N v_{ij}^2 + \Psi_j$

So $b_{jj} = \frac{1}{2}(d_{jj}^2 + d_{jj}^2 - d_{jj}^2) = \mathbf{B} = \mathbf{X}\mathbf{X}^T$ is p.s.d. :

$\mathbf{B} = \mathbf{C}\mathbf{D}^2\mathbf{C}^T = (\mathbf{C}\mathbf{D}^{1/2})(\mathbf{C}\mathbf{D}^{1/2})^T$

Ignore small eigenvalues, let \mathbf{c}_j be the k eigenvectors chosen with eigenvalues λ_j , the new dimensions: $\mathbf{z}_j^{(i)} = \sqrt{\lambda_j} \mathbf{c}_j^T \mathbf{x}^{(i)}$

LDA Sample mean after projection:

$\mathbf{m}_1 = \mathbf{w}^T \mathbf{m}_1, \mathbf{m}_2 = \mathbf{w}^T \mathbf{m}_2$

Between class scatter: $(\mathbf{m}_1 - \mathbf{m}_2) = \mathbf{w}^T \mathbf{S}_B \mathbf{w}$, $\mathbf{S}_B = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T$

Within Class Scatter: $\mathbf{S}_1^2 = \mathbf{w}^T \mathbf{S}_1 \mathbf{w}$, $\mathbf{S}_1 = \sum_i (\mathbf{x}^{(i)} -$

log likelihood: $\mathcal{L} \equiv \log L(\theta|\mathcal{X}) = \sum_{i=1}^N \log p(\mathbf{x}^{(i)}|\theta)$

Max Likelihood estimate: $\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta|\mathcal{X})$

Bernoulli $\mathbf{x} \in \{0, 1\}$, $P(x = 1)$ for $p(C_1)$,
 $P(\mathbf{x}|\mathbf{p}_0) = p_0^{\sum x^{(i)}} (1 - p_0)^{N - \sum x^{(i)}}$

Log Likelihood:

$\mathcal{L}(\mathbf{p}_0|\mathcal{X}) = \sum_{i=1}^N [x^{(i)} \log p_0 + (1 - x^{(i)}) \log(1 - p_0)]$

ML estimation: $\hat{p}_0 = \frac{1}{N} \sum_{i=1}^N x^{(i)}$

Multinomial Ran var. \mathbf{x} with $K \geq 2$ possible value

Indicator var: $x_i = 1$ if outcome is state i , 0 if not

$P(\mathbf{x}|\mathbf{p}_1, \dots, \mathbf{p}_K | p_1, \dots, p_K) = \prod_{i=1}^N p_i^{x_i}, \sum_{i=1}^K p_i = 1$

Log likelihood: $\mathcal{L}(\mathbf{p}_0|\mathcal{X}) = \sum_{i=1}^N \sum_{k=1}^K x_{ik}^{(i)} \log p_{ik}$

ML estimate: $\hat{p}_k = \frac{1}{N} \sum_{i=1}^N x_{ik}^{(i)}$

Normal pdf: $p(\mathbf{x}|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp[-\frac{(\mathbf{x}-\mu)^2}{2\sigma^2}]$

$\mathcal{L}(\mu, \sigma|\mathcal{X}) = -\frac{N}{2} \log(2\pi) - N \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^N (\mathbf{x}^{(i)} - \mu)^2$

ML Estimates: $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)}$

$\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}^{(i)} - \hat{\mu})^2$

Multivariable Normal $\mathcal{N}(\mu, \Sigma)$, μ mean vec, Σ covariance matrix

$p(\mathbf{x}|\mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp[-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)]$

$\mathcal{L}(\mu, \Sigma|\mathcal{X}) = -\frac{Nd}{2} \log(2\pi) - \frac{N}{2} \log |\Sigma| - \frac{1}{2} \sum_{i=1}^N (\mathbf{x}^{(i)} - \mu)^T \Sigma^{-1}(\mathbf{x}^{(i)} - \mu)$

ML estimates: $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)}$

$\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}^{(i)} - \hat{\mu})(\mathbf{x}^{(i)} - \hat{\mu})^T$

Bias and Variance Bias: $b_f(d) = E[f] - \theta$

Variance: $E[d - E[d]^2]$ (d is estimator of param θ)

Mean Squared error:

$r(d, \theta) = E[(d - \theta)^2] = \text{bias}^2 + \text{variance}^2$

Bayesian Estimation Treat θ as a Ran Var. Prior $p(\theta)$

Posterior:

$p(\theta|\mathcal{X}) = \frac{p(\mathcal{X}|\theta)p(\theta)}{p(\mathcal{X})} = \frac{p(\mathcal{X}|\theta)p(\theta)}{\int p(\mathcal{X}|\theta')p(\theta')d\theta'}$

Estimation of density at \mathbf{x} :

$p(\mathbf{x}|\mathcal{X}) = \int p(\mathbf{x}|\theta)p(\theta|\mathcal{X})d\theta$

Regression $y = g(\mathbf{x}|\theta)$: $y = \int g(\mathbf{x}|\theta)p(\theta|\mathcal{X})d\theta$

Computational Considerations

Max a posteriori (MAP): $\theta_{MAP} = \arg \max_{\theta} p(\theta|\mathcal{X})$,

$p(\mathbf{x}|\mathcal{X}) \approx p(\mathbf{x}|\theta_{MAP})$, $y \approx y_{MAP} = g(\mathbf{x}|\theta_{MAP})$

ML estimation: $\theta_{ML} = \arg \max_{\theta} p(\theta|\mathcal{X})$

Bayes' estimation - expectation w.r.t. posterior density:

$\theta_{Bayes} = E[\theta|\mathcal{X}] = \int \theta p(\theta|\mathcal{X})d\theta$

Example: Bayesian estimation with known μ , σ and σ_0

$x^{(i)} \sim \mathcal{N}(\theta, \sigma_0^2)$, $\theta \sim \mathcal{N}(\mu, \sigma^2)$

MLE: $\theta_{ML} = \frac{1}{N} \sum_{i=1}^N x^{(i)} = m$

$\theta_{Map} = \theta_{Bayes} = E(\theta|\mathcal{X}) = \frac{N/\sigma^2}{N/\sigma^2 + 1/\sigma^2} m + \frac{1/\sigma^2}{N/\sigma^2 + 1/\sigma^2} \mu$

$\mathbf{m}_1(\mathbf{x}^{(i)} - \mathbf{m}_1)^T y^{(i)}$, similarly $\mathbf{S}_2 = \sum_i (\mathbf{x}^{(i)} - \mathbf{m}_2)(\mathbf{x}^{(i)} - \mathbf{m}_2)^T (1 - y^{(i)})$, so

$\mathbf{S}_1^2 + \mathbf{S}_2^2 = \mathbf{w}^T \mathbf{S}_B \mathbf{w}$

$\mathbf{S}_B = \mathbf{S}_1 + \mathbf{S}_2$

Fisher's LD $J(\mathbf{w}) = \frac{(\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}}{\sigma_1^2 + \sigma_2^2} = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}$, take derivative of J w.r.t. \mathbf{w} setting it to 0: $\mathbf{S}_B \mathbf{w} = \lambda \mathbf{S}_B \mathbf{w}$,

or $\mathbf{S}_B^{-1} \mathbf{S}_B \mathbf{w} = \lambda \mathbf{w}$ (Eigen equation)

$K > 2$ Within-class scatter $\mathbf{S}_1 = \sum_i y_i^{(i)} (\mathbf{x}^{(i)} - \mathbf{m}_1)(\mathbf{x}^{(i)} - \mathbf{m}_1)^T$,

$G^0 = 1$ if $\mathbf{x}^{(i)} \in C_i$

Total class scatter: $\mathbf{S}_B = \sum_{i=1}^K \mathbf{S}_i$

Between Class Scatter: $\mathbf{S}_B = \sum_{i=1}^K N_i (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^T$

Optimal is \mathbf{W} that max: $J(\mathbf{W}) = \frac{\mathbf{W}^T (\mathbf{W}^T \mathbf{S}_B \mathbf{W})}{\mathbf{W}^T (\mathbf{W}^T \mathbf{S}_B \mathbf{W})}$ Corresponds to eigenvectors of $\mathbf{S}_B^{-1} \mathbf{S}_B$

5 Clustering

Mixture Densities $p(\mathbf{x}) = \sum_{i=1}^K p(\mathbf{x}|\mathcal{G}_i)P(\mathcal{G}_i)$ \mathcal{G}_i mixture components, $p(\mathbf{x}|\mathcal{G}_i)$ component densities, $P(\mathcal{G}_i)$ mixing proportions

Gaussian mixture: $p(\mathbf{x}|\mathcal{G}_i) = \mathcal{N}(\mu_i, \Sigma_i)$,

parameters $\Phi = \{p(\mathcal{G}_i), \mu_i, \Sigma_i\}_{i=1}^K$

k-Means Initialize $m_i, i = 1, \dots, k$ (random)

Repeat

For all $\mathbf{x} \in \mathcal{X}$

$b_i^{(t)} = \begin{cases} 1 & \text{if } i = \arg \min_j \|\mathbf{x}^{(i)} - \mathbf{m}_j\| \\ 0 & \text{otherwise} \end{cases}$

For all $\mathbf{m}_i, i = 1, \dots, k$

$\mathbf{m}_i = \frac{\sum_j b_j^{(t)} \mathbf{x}^{(j)}}{\sum_j b_j^{(t)}}$

Until \mathbf{m}_i converges

EM Algorithm

Log likelihood:

$\mathcal{L}(\Phi|\mathcal{X}) = \log \prod_i p(\mathbf{x}^{(i)}|\Phi) = \sum_i \log \sum_j p(\mathbf{x}^{(i)}|\mathcal{G}_j)P(\mathcal{G}_j)$

where Φ include $P(\mathcal{G}_i)$ and sufficient statistic Θ_i of $p(\mathbf{x}^{(i)}|\mathcal{G}_i)$

Complete-data likelihood: $\mathcal{L}_C(\Phi|\mathcal{X}, \mathcal{Z})$, \mathcal{Z} hidden variables (not observable)

Expectation of complete data likelihood given \mathcal{X} and the current (iteration t) parameter val $\Phi^{(t)}$ (E-step)

$Q(\Phi|\Phi^{(t)}) = E[\mathcal{L}_C(\Phi|\mathcal{X}, \mathcal{Z})|\mathcal{X}, \Phi^{(t)}]$

$= \sum_j p(\mathcal{Z}|\mathcal{X}, \Phi^{(t)}) \log p(\mathcal{Z}, \Phi)$

m-step: $\Phi^{(t+1)} = \arg \max_{\Phi} Q(\Phi, \Phi^{(t)})$

EM in Gaussian Mixtures Indicator $\mathbf{z}^{(i)}$, $z_i^{(i)} = 1$ if $\mathbf{x}^{(i)}$ belongs to cluster \mathcal{G}_i

Gaussian Component densities: $p_i(\mathbf{x}^{(i)}) = p(\mathbf{x}|\mathcal{G}_i) = \mathcal{N}(\mu_i, \Sigma_i)$

Prior $P(\mathcal{G}_i) = \pi_i$, so: $p(\mathbf{x}^{(i)}) = \prod_{i=1}^k \pi_i^{z_i^{(i)}}$

Classification with Discriminant Functions, Gaussian density for each class: $p(\mathbf{x}|C_i) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp[-\frac{(\mathbf{x}-\mu_i)^2}{2\sigma_i^2}]$

Discriminant functions: $g_i(\mathbf{x}) = \log [p(\mathbf{x}|C_i)P(C_i)] = -\frac{1}{2} \log 2\pi - \log \sigma_i - \frac{(\mathbf{x}-\mu_i)^2}{2\sigma_i^2} + \log P(C_i)$

Sample $\mathcal{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$, $g_i^{(i)} = 1$ if $\mathbf{x}^{(i)} \in C_i$

ML Estimates: $\hat{P}(C_i) = \frac{1}{N} \sum_{i=1}^N g_i^{(i)}$

$m_i = \frac{\sum_{i=1}^N g_i^{(i)} \mathbf{x}^{(i)}}{\sum_{i=1}^N g_i^{(i)}}$, $s_i^2 = \frac{\sum_{i=1}^N (g_i^{(i)} - m_i)^2 \mathbf{x}^{(i)^2}}{\sum_{i=1}^N g_i^{(i)}}$

Discriminant Functions:

$g_i(\mathbf{x}) = -\log s_i - \frac{(\mathbf{x} - m_i)^2}{2s_i^2} + \log \hat{P}(C_i)$

Additive Parametric Model Functional relationship in additive form: $r = f(\mathbf{x}) + \epsilon$

Parametric modeling: $f(\mathbf{x}) \approx g(\mathbf{x}|\theta)$, $\epsilon \sim \mathcal{N}(0, \sigma^2)$

Conditional probability of output given input:

$p(r|\mathbf{x}) \sim \mathcal{N}(g(\mathbf{x}|\theta), \sigma^2)$

Log likelihood given: $\mathcal{X} = \{(\mathbf{x}^{(i)}, r^{(i)})\}$:

$\mathcal{L}(\theta|\mathcal{X}) = \log \prod_{i=1}^N p(r^{(i)}|\mathbf{x}^{(i)}) = \frac{1}{2\sigma^2} \sum_{i=1}^N [r^{(i)} - g(\mathbf{x}^{(i)}|\theta)]$

Geometry Interpretation Express any point as $\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{x}_p - \mathbf{x}_o}{\|\mathbf{x}_p - \mathbf{x}_o\|}$, \mathbf{x}_p projection of \mathbf{x} onto hyperplane. r distance from \mathbf{x} to hyper plane, we have $r = \frac{|\mathbf{w}^T \mathbf{x} + w_0|}{\|\mathbf{w}\|}$

Multi classes K discriminant $g_i(\mathbf{x}|\mathbf{w}_i, w_{i0}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$, linear separable $g_i(\mathbf{x}|\mathbf{w}_i, w_{i0}) > 0$ if $\mathbf{x} \in C_i$, choose C_i if $g_i(\mathbf{x}) = \max_{j=1}^K g_j(\mathbf{x})$

Pairwise Separation Discriminant function for class i and j :

$$g_{ij}(\mathbf{x}|\mathbf{w}_{ij}, w_{ij0}) = \mathbf{w}_{ij}^T \mathbf{x} + w_{ij0} = \begin{cases} > 0 & \text{if } \mathbf{x} \in C_i \\ \leq 0 & \text{if } \mathbf{x} \in C_j \\ \text{don't care} & \text{if } \mathbf{x} \in C_k, k \neq i, k \neq j \end{cases}$$

Logistic Discrimination

Two Classes Assume that the log likelihood ratio is linear:

$$\log \frac{p(\mathbf{x}|C_1)}{p(\mathbf{x}|C_2)} = \mathbf{w}^T \mathbf{x} + w_0$$

$$\log(P(C_1|\mathbf{x})) = \log \frac{p(C_1|\mathbf{x})}{1 - p(C_1|\mathbf{x})}$$

$$= \log \frac{p(\mathbf{x}|C_1)}{p(\mathbf{x}|C_2)} + \log \frac{p(C_1)}{p(C_2)}$$

$$= \mathbf{w}^T \mathbf{x} + w_0$$

where $w_0 = w_0^0 + \log \frac{p(C_1)}{p(C_2)}$ Rearranging terms:

$$y = \text{sigmoid}(\mathbf{w}^T \mathbf{x} + w_0)$$

$$= \frac{1}{1 + \exp[-(\mathbf{w}^T \mathbf{x} + w_0)]}$$

As our estimator of $P(C_1|\mathbf{x})$

Gradient Decent In the discriminant-based approach, the parameters are those of the discriminants, and they are *optimized to minimize the classification error*

Error w denotes the set of parameters and $E(\mathbf{w}|\mathcal{X})$ is the error parameters \mathbf{w} on the given training set \mathcal{X} , we look for:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} E(\mathbf{w}|\mathcal{X})$$

No analytical solution

Gradient Vector When $E(\mathbf{w})$ is a differentiable function of a vector of variables, we have the gradient vector composed of the partial derivatives:

$$\nabla_{\mathbf{w}} E = \left[\frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_d} \right]^T$$

Gradient Decent starts from a *random* \mathbf{w} , at each step, update \mathbf{w} in a *opposite direction* of the gradient:

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}, \forall i$$

$$w_i = w_i + \Delta w_i$$

η is *step size*, or *learning factor*

When we get to minimum, the derivative is 0 and the procedure terminates.

This indicates that the procedure finds the nearest minimum that can be *local minimum*. There is no guarantee of finds the nearest minimum that can be a local minimum

Learning parameters Given a sample of two classes, $\mathcal{X} = \mathbf{x}^{(i)}, \mathbf{r}^{(i)}$, where $\mathbf{r}^{(i)} = 1$ if $\mathbf{x} \in C_1$

We assume $\mathbf{r}^{(i)}$ given $\mathbf{x}^{(i)}$ is Bernoulli with probability $y^{(i)} = p(C_1|\mathbf{x}^{(i)})$:

$$\mathbf{r}^{(i)}|\mathbf{x}^{(i)} \sim \text{Bernoulli}(y^{(i)})$$

Note that in this discriminant-based approach, we model directly $\mathbf{r}|\mathbf{x}$ The sample likelihood is:

$$L(\mathbf{w}, \mathbf{w}_0|\mathcal{X}) = \prod_i (y^{(i)})^{r^{(i)}} (1 - y^{(i)})^{1-r^{(i)}}$$

We can always turn it in an error function to minimize: $E = -\log L$ So we have *cross-entropy*:

$$E(\mathbf{w}, \mathbf{w}_0|\mathcal{X}) = -\sum_i \mathbf{r}^{(i)} \log y^{(i)} + (1 - \mathbf{r}^{(i)}) \log(1 - y^{(i)})$$

We use gradient descent to minimize cross-entropy If $y = \text{sigmoid}(a)$ is $\frac{1}{1 + \exp(-a)}$, its derivative is given as:

$$\frac{dy}{da} = y(1 - y)$$

and we get the following update equations:

$$\Delta w_j = -\eta \frac{\partial E}{\partial w_j} = \eta \sum_i \mathbf{r}^{(i)} \left(\frac{1}{y^{(i)}} - \frac{1 - \mathbf{r}^{(i)}}{1 - y^{(i)}} \right) x_j^{(i)}$$

$$= \eta \sum_i (\mathbf{r}^{(i)} - y^{(i)}) x_j^{(i)}, i = 1, \dots, d$$

$$\Delta w_0 = -\eta \frac{\partial E}{\partial w_0} = \eta \sum_i (\mathbf{r}^{(i)} - y^{(i)})$$

Multiple Classes Take one of the classes C_k , as reference class and assume that:

$$\log \frac{p(\mathbf{x}|C_k)}{p(\mathbf{x}|C_1)} = \mathbf{w}_k^T \mathbf{x} + w_{k0}, i = 1, \dots, K - 1$$

Then we have:

$$\frac{P(C_k|\mathbf{x})}{P(C_1|\mathbf{x})} = \exp(\mathbf{w}_k^T \mathbf{x} + w_{k0})$$

With $w_{k0} = w_{k0}^0 + \log \frac{p(C_k)}{p(C_1)}$

Generalization of sigmoid Summing over i we can deduce:

$$\frac{P(C_k|\mathbf{x})}{P(C_1|\mathbf{x})} = \frac{\exp(\mathbf{w}_k^T \mathbf{x} + w_{k0})}{1 + \sum_{i=1}^{K-1} \exp(\mathbf{w}_i^T \mathbf{x} + w_{i0})}, k = 1, \dots, K - 1$$

Softmax Treat all classes uniformly

$$y_i = \frac{P(C_i|\mathbf{x})}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \mathbf{x} + w_{j0})}, i = 1, \dots, K$$

Learning $\frac{\partial w_i}{\partial w_j} = y_i(\delta_j - y_j)$, $\Delta \mathbf{w}_j = \eta \sum_i (\mathbf{r}_j^{(i)} - y_j^{(i)}) \mathbf{x}^{(i)}$, $\Delta w_{j0} = \eta \sum_i (\mathbf{r}_j^{(i)} - y_j^{(i)})$

Regression for two-class Classification $\mathbf{r}^{(i)} = y^{(i)} + \epsilon$, $\mathbf{r}^{(i)} \in \{0, 1\}$, $\epsilon \sim \mathcal{N}(0, \sigma^2)$, $y^{(i)} = \text{sigmoid}(\mathbf{w}^T \mathbf{x}^{(i)} + w_0)$,

$$\text{Likelihood: } L(\mathbf{w}, w_0|\mathcal{X}) = \prod_i \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(\mathbf{r}^{(i)} - y^{(i)})^2}{2\sigma^2} \right],$$

$$\text{Error: } E(\{\mathbf{w}_i, w_{i0}\}, |\mathcal{X}) = \frac{1}{2} \sum_i (\mathbf{r}^{(i)} - y^{(i)})^2,$$

$$\text{Learning } \Delta \mathbf{w} = \eta \sum_i (\mathbf{r}^{(i)} - y^{(i)}) y^{(i)} (1 - y^{(i)}) \mathbf{x}^{(i)},$$

$$\Delta w_0 = \eta \sum_i (\mathbf{r}^{(i)} - y^{(i)}) y^{(i)} (1 - y^{(i)})$$

11 Performance Evaluation and Comparison

K-Fold Cross Validation The data set \mathcal{X} is randomly partitioned into K equal-sized subsets X_i . **Stratification** if the class distribution different subsets are kept roughly the same. Use X_1, \dots, X_K as validation sets, and the remaining as training respectively. If N is small, K should be large to allow large enough training sets. Leave one out: one instance if left out as validation and $N - 1$ for training.

5 x 2 Cross Validation For each fold i , Split into two equal-sized parts, $X_1^{(i)}, X_2^{(i)}$, 10 training/validation set pairs: $T_1 = X_1^{(1)}, V_1 = X_2^{(1)}, T_2 = X_1^{(2)}, V_2 = X_2^{(2)}, T_3 = X_2^{(1)}, V_3 = X_1^{(1)}, T_4 = X_2^{(2)}, V_4 = X_1^{(2)}, \dots$

Bootstrapping Generates new samples, each of size N , by drawing randomly with replacement. Multiple bootstrap samples are used to max the change that the system is trained on all instances.

Error Measure TP: True positive, FP: False positive, FN: False Negative, TN: True Negative, Error rate: $\frac{FP + FN}{N}$

Receiver Operating characteristics curve Hit rate: $\frac{TP}{TP + FN}$, False alarm rate: $\frac{FP}{FP + TN}$, Area under the curve is often used.

Point vs Interval Estimator Point specifies a value for θ , interval specifies an interval within which θ lines with a certain degree of confidence.

Confidence Interval Define z_{α} ,

Two sided: $P(Z > z_{\alpha}) = \alpha$, then for level of confidence $1 - \alpha$,

$$P(-z_{\alpha/2} < \sqrt{N} \frac{\bar{m} - \mu}{\sigma} < z_{\alpha/2}) = 1 - \alpha$$

$$\text{One-sided } P(m - z_{\alpha} \sigma \sqrt{\frac{m}{N}} < \mu) = 1 - \alpha$$

t Distribution σ^2 replaced by sample var $S^2 = \frac{\sum_i (x_i^{(i)} - \bar{m})^2}{N - 1}$, $\sqrt{N}(\bar{m} - \mu)/S$ follows a t distribution with $N - 1$ degree of freedom: $\sqrt{N} \frac{\bar{m} - \mu}{S} \sim t_{N-1}$ For $\alpha \in (0, 1/2)$,

$$P(m - t_{\alpha/2, N-1} \cdot \frac{S}{\sqrt{N}} < \mu < m + t_{\alpha/2, N-1} \cdot \frac{S}{\sqrt{N}}) = 1 - \alpha$$

Hypothesis Testing Test some hypothesis concerning the parameters.

1. Define a statistics obey a certain distribution
2. Random sample is consistent with the hypothesis under consideration, accept(not reject)
3. Otherwise reject

Given normal $N(\mu, \sigma^2)$, σ known, μ unknown. **Null Hypothesis:** $H_0: \mu = \mu_0$, $H_1: \mu \neq \mu_0$ Accept H_0 if the sample mean is not too far from μ_0 . Accept H_0 with level of significance α if μ_0 lies in the $100(1 - \alpha)$.

$$\sqrt{N} \frac{\bar{m} - \mu_0}{\sigma} \in (-z_{\alpha/2}, z_{\alpha/2})$$

$$K > 2 \text{ classes } \mathbf{r}^{(i)} = \mathbf{y}^{(i)} + \epsilon,$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2 K), y^{(i)} = \frac{1}{1 + \exp[-(\mathbf{w}^T \mathbf{x}^{(i)} + w_0)]},$$

$$\text{Likelihood: } L(\{\mathbf{w}_i, w_{i0}\}, |\mathcal{X}) = \prod_i \prod_{j \in \{1, \dots, K\} \setminus \{i\}} \exp \left[-\frac{(\mathbf{r}^{(i)} - y_j^{(i)})^2}{2\sigma^2} \right]$$

$$\text{Error Func: } E(\{\mathbf{w}_i, w_{i0}\}, |\mathcal{X}) = \frac{1}{2} \sum_i \|\mathbf{r}^{(i)} - \mathbf{y}^{(i)}\|^2$$

$$\text{Learning } \Delta \mathbf{w}_i = \eta \sum_i (\mathbf{r}^{(i)} - y_i^{(i)}) y_i^{(i)} (1 - y_i^{(i)}) \mathbf{x}^{(i)},$$

$$\Delta w_0 = \eta \sum_i (\mathbf{r}^{(i)} - y_i^{(i)}) y_i^{(i)} (1 - y_i^{(i)})$$

9 Multilayer Perceptrons

Perceptron The output y is a **weighted sum** of the input $\mathbf{x} = (x_0, x_1, \dots, x_d)^T$:

$$y = \sum_{j=1}^d w_j x_j + w_0 = \mathbf{w}^T \mathbf{x}$$

where x_0 is a special *bias unit* with $x_0 = 1$ and $\mathbf{w} = (w_0, w_1, \dots, w_d)^T$ are called **connection weight** or **synaptic weights**

to implement a linear discriminant function, we need threshold function:

$$s(a) = \begin{cases} 1 & \text{if } a > 0 \\ 0 & \text{otherwise} \end{cases}$$

to define the following decision rule:

$$\text{Choose } \begin{cases} C_1 & \text{if } s(\mathbf{w}^T \mathbf{x}) = 1 \\ C - 2 & \text{otherwise} \end{cases}$$

Use **sigmoid** instead of threshold to gain differentiability:

$$y = \text{sigmoid}(\mathbf{w}^T \mathbf{x}), \text{ where } \text{sigmoid}(a) = \frac{1}{1 + \exp(-a)}$$

The output may be interpreted as the *posterior probability* that the input \mathbf{x} belongs to C_1

$K > 2$ **Outputs** K perceptrons, each with a weight vector \mathbf{w}_i ,

$$y_i = \sum_{j=1}^d w_{ij} x_j + w_{i0} = \mathbf{w}_i^T \mathbf{x}, \text{ or } \mathbf{y} = \mathbf{W} \mathbf{x}$$

where w_{ij} is the weight from input x_j to output y_i and each row of the $K \times (d + 1)$ matrix \mathbf{W} is the weight vector of one perceptron.

Choose C_i if $y_i = \max_k y_k$

Posterior probability, use softmax to define y_i as:

$$y_i = \frac{\exp(\mathbf{w}_i^T \mathbf{x})}{\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x})}$$

Stochastic Gradient Decent Gradient descent for online learning, for regression, the error on a single instance

$$E^{(i)}(\mathbf{w}|\mathbf{x}^{(i)}, \mathbf{r}^{(i)}) = \frac{1}{2} (\mathbf{r}^{(i)} - y^{(i)})^2 = \frac{1}{2} (\mathbf{r}^{(i)} - (\mathbf{w}^T \mathbf{x}^{(i)}))^2$$

, gives the online update rule:

$$\Delta \mathbf{w}_j^{(i)} = \eta (\mathbf{r}^{(i)} - y^{(i)}) x_j^{(i)}, \eta \text{ is step size.}$$

For Binary classification:

$$\text{Likelihood: } L = (y^{(i)})^{r^{(i)}} (1 - y^{(i)})^{1-r^{(i)}}$$

$$\text{Cross Entropy: } E(\mathbf{w}|\mathbf{x}^{(i)}, \mathbf{r}^{(i)}) = -\log L = -r^{(i)} \log y^{(i)} - (1 - r^{(i)}) \log(1 - y^{(i)})$$

$$\text{Online update rule: } \Delta \mathbf{w}_j^{(i)} = \eta (\mathbf{r}^{(i)} - y^{(i)}) x_j^{(i)}$$

Error Type I, rejected when it is correct, α defines how much type I can tolerate. Type II accepted when incorrect.

One Sided $H_0: \mu \leq \mu_0$, $H_1: \mu > \mu_0$, accept $\frac{\sqrt{N}(\bar{m} - \mu_0)}{\sigma} \in (-\infty, z_{\alpha})$

t-test If σ^2 is not known, $H_0: \mu = \mu_0$, $H_1: \mu \neq \mu_0$, accept at α if

$$\sqrt{\frac{N}{N-1}} \frac{\bar{m} - \mu_0}{S} \in (-t_{\alpha, N-1}, t_{\alpha, N-1})$$

Binomial Test Single test/validation. Classifier is trained on a training set T and tested on validation set V . p be the probability that the classifier makes a misclassification error. Define $x^{(i)}$ as a Bernoulli variable to denote the correctness, $x^{(i)} = 1$ with probability p and $x^{(i)} = 0$ with probability $1 - p$, point estimation: $\hat{p} = \frac{\sum_i x^{(i)}}{N} = \frac{\sum_i x^{(i)}}{N}$

Hypothesis test: $H_0: p \leq p_0$ vs. $H_1: p > p_0$, $X = \sum_{i=1}^N x^{(i)}$ denote the number of errors on V ,

$$P(X = j) = \binom{N}{j} p^j (1 - p)^{N-j}$$

Under the null hypothesis $p \leq p_0$, so the probability that there are e errors or less is:

$$P(X \leq e) = \sum_{j=1}^e \binom{N}{j} p_0^j (1 - p_0)^{N-j}$$

Binomial test: accept H_0 if $P(X \leq e) < 1 - \alpha$

Paired t Test Multiple training/validation set pairs: run the algorithm K times on K/T set pairs, we get K error probabilities, $p_i, i = 1, \dots, K$ on K validation sets.

Paired t test to determine whether to accept the null hypothesis H_0 that the classifier has error probability p_0 or less at significance level α .

$$x_1^{(i)} = \begin{cases} 1 & \text{if classifier trained on } V_i \text{ makes an error on instance } i \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Test statistic: } \sqrt{K} \frac{(\bar{m} - p_0)}{\sigma} \sim t_{K-1}, \text{ where } p_0 = \frac{\sum_{i=1}^K x_1^{(i)}}{K}, m = \frac{\sum_{i=1}^K x_1}{K}, S^2 = \frac{\sum_{i=1}^K (x_1 - m)^2}{K-1}$$

K-Fold Cross-Validated Paired t Test Given two classification algorithms and a data set, we want to compare and test whether the two algorithms construct classifiers that have the same expected error rate on a new instance.

K-fold cross validation is used to obtain K training/validation set pairs $\{(\mathcal{T}_i, \mathcal{V}_i)\}_{i=1}^K$, run two algorithms, error probabilities p_i and p_i' , define $p_i = p_i' - p_i$, p_i with mean μ

$H_0 = \mu = 1, H_1 = \mu \neq 0$. Test statistic: $\sqrt{K} \frac{(\bar{m} - 0)}{\sigma} \sim t_{K-1}$, where $m = \frac{\sum_{i=1}^K p_i}{K}, S^2 = \frac{\sum_{i=1}^K (p_i - m)^2}{K-1}$, accept H_0 at significance level α if $\sqrt{K} m / S \in -t_{\alpha/2, K-1}, t_{\alpha/2, K-1}$

12 Ensemble Learning

Combining Base Learners Multicexpert combination method: base learners work in parallel, give decision and combined to give a final.

Multi-stage combinations: Base learners work serially, sorted increasing complexity, complex is used when simple is not confident.

$$K > 2 \text{ } y_i^{(i)} = \frac{\exp(\mathbf{w}_i^T \mathbf{x}^{(i)})}{\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x}^{(i)})}$$

$$\text{Likelihood: } L = \prod_i (y_i^{(i)})^{r^{(i)}}$$

$$\text{Cross Entropy: } E^{(i)}(\{\mathbf{w}_i\}, \mathbf{x}^{(i)}, \mathbf{r}^{(i)}) = -\sum_j r_j^{(i)} \log y_j^{(i)}$$

$$\text{Online update rule: } \Delta \mathbf{w}_j^{(i)} = \eta (\mathbf{r}^{(i)} - y_j^{(i)}) x_j^{(i)}$$

Multilayer Perceptron MLP has a hidden layer between the input and output. Input to hidden: $z_h = \text{sigmoid}(\mathbf{w}_h^T \mathbf{x}) = \frac{1}{1 + \exp[-(\sum_{j=1}^d w_{hj} x_j + w_{h0})]}$

Hidden to output: $y_i = \mathbf{v}_i^T \mathbf{z} = \sum_{h=1}^H v_{ih} z_h + v_{i0}$

Backward, hidden-to-output weight: treating hidden unit as input

input-to-hidden, chain rule: $\frac{\partial E}{\partial w_{hj}} = \frac{\partial E}{\partial y_h} \frac{\partial y_h}{\partial z_h} \frac{\partial z_h}{\partial w_{hj}}$

MLP for Nonlinear Regression(Multi output) outputs: $y_i^{(i)} = \sum_{h=1}^H v_{ih} z_h^{(i)} + v_{i0}$, $z_h^{(i)} = \text{sigmoid}(\mathbf{w}_h^T \mathbf{x}^{(i)})$

$$\text{Error function: } E(\mathbf{W}, \mathbf{v}|\mathcal{X}) = \frac{1}{2} \sum_i \sum_j (\mathbf{r}_j^{(i)} - y_j^{(i)})^2$$

$$\text{Update rule for second layer: } \Delta v_{ih} = \eta \sum_j (\mathbf{r}_j^{(i)} - y_j^{(i)}) z_h^{(i)}$$

$$\text{Update rule for first layer: } \Delta w_{hj} = -\eta \frac{\partial E}{\partial w_{hj}} = -\eta \sum_i (\mathbf{r}_i^{(i)} - y_i^{(i)}) v_{ih} z_h^{(i)} (1 - z_h^{(i)}) x_j^{(i)}$$

MLP for NonLinear Multi-class Discrimination

$$\text{Outputs: } y_i^{(i)} = \frac{\exp(\mathbf{w}_i^T \mathbf{x}^{(i)})}{\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x}^{(i)})}, \mathbf{a}_i^{(i)} = \sum_{k=1}^K v_{ik} z_k^{(i)} + v_{i0}$$

$$\text{Error Function: } E(\mathbf{W}, \mathbf{V}|\mathcal{X}) = -\sum_i \sum_j r_j^{(i)} \log y_j^{(i)}, \text{ update rules are the same as regression.}$$

10 Support Vector Machine

Optimal Separating Hyperplane results from *statistical learning* theory showing that the separating hyperplane with the largest margin generalizes best.

Hard-margin case points are **linearly separable** with proper scaling of \mathbf{w} and w_0 , the points closest to the hyper plane satisfy $|\mathbf{w}^T \mathbf{x} + w_0| = 1 \Rightarrow$ **canonical separating hyperplane**.

The one max the margin: **canonical optimal separating hyperplane**

$\mathbf{x}^{(1)}$ and $\mathbf{x}^{(2)}$ be two closest point on each side: