

Probability

- Bernoulli (π) $E = \pi$ - Uniform $U(a,b)$ $E = \frac{a+b}{2}$ $Var = \frac{(b-a)^2}{12}$
 $P(Y=y) = \begin{cases} \pi & y=1 \\ 1-\pi & y=0 \end{cases}$ $Var = \pi(1-\pi)$ $P(Y=y) = \begin{cases} \frac{1}{b-a} & \text{if } y \in [a,b] \\ 0 & \text{otherwise} \end{cases}$
- Binomial (n, π) $E = n\pi$, $Var = n\pi(1-\pi)$ - normal distribution $N(\mu, \sigma^2)$ $E = \mu$
 $P(Y=y) = \binom{n}{y} \pi^y (1-\pi)^{n-y}$ $P(Y=y) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2}(\frac{y-\mu}{\sigma})^2}$ $Var = \sigma^2$
- Poisson (λ) $E = \lambda$, $Var = \lambda$ - standard normal distribution $N(0,1)$
 $P(Y=y) = \frac{e^{-\lambda} \lambda^y}{y!}$ $P(Y=y) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}y^2}$

conditional probability

Bayes rule
 $P(Y|X) = \frac{P(X|Y) \cdot P(Y)}{P(X)}$
 $P(Y|X) = \frac{P(X|Y) \cdot P(Y)}{\sum_y P(X|y) \cdot P(y)}$

example:
 $P(HIV|pos) = \frac{P(pos|HIV) \cdot P(HIV)}{P(pos)} = \frac{P(pos|HIV) \cdot P(HIV)}{P(pos|HIV) \cdot P(HIV) + P(pos|not HIV) \cdot P(not HIV)}$

- independence
 two events are independent if $P(X|Y) = P(X)$
 $P(X,Y) = P(X) \cdot P(Y|X) = P(X) \cdot P(Y)$
- pdf vs cdf probability density function cumulative distribution function
 $F_X(y) = P(X \leq y)$ ordered.
 $\int_{-\infty}^y f_X(x) dx$ continuous
- For random variable $Y \sim N(\mu, \sigma^2)$
 $P(Y \leq y) = P\left\{\frac{Y-\mu}{\sigma} \leq \frac{y-\mu}{\sigma}\right\} = P\left\{Z \leq \frac{y-\mu}{\sigma}\right\}$

Bivariate Normal distribution (Y_1, Y_2)
 requires the parameters $\mu_1, \sigma_1, \mu_2, \sigma_2, \rho$
 $\rho_{12} = \frac{Cov(Y_1, Y_2)}{\sigma_1 \sigma_2} = \frac{E[(Y_1 - \mu_1)(Y_2 - \mu_2)]}{\sigma_1 \sigma_2} = \frac{E(Y_1 Y_2) - E(Y_1)E(Y_2)}{\sigma_1 \sigma_2}$
 coefficient of correlation

- Uncorrelated \neq independent
 independent \rightarrow uncorrelated
 uncorrelated \nrightarrow independent.
- $E(Y_1 Y_2) = (\mu_1 - \rho_{12} \frac{\sigma_1}{\sigma_2}) \mu_2 + \rho_{12} \frac{\sigma_1}{\sigma_2} \mu_2 = \mu_1 \mu_2 + \rho_{12} \frac{\sigma_1}{\sigma_2} \mu_2$
 $Var(Y_1 Y_2) = \sigma_1^2 (1 - \rho_{12}^2)$

Information theory

- Entropy \rightarrow measure of its uncertainty
 $H(y) = - \sum_{c=1}^C P(y=c) \log_2 P(y=c)$
- Kullback-Leibler divergence
 - dissimilarity of two prob. distributions p and q
 $KL(p,q) = \sum_{c=1}^C P_c \log_2 \frac{P_c}{Q_c} = \sum_{c=1}^C P_c \log_2 p_c - \sum_{c=1}^C P_c \log_2 q_c$
 $= -H(p) + H(p,q)$
 $= -\text{entropy} + \text{cross-entropy}$
- Mutual information
 - MI = 0 iff the variables are independent
 $I(x,y) = \sum_{x,y} P(x,y) \log_2 \frac{P(x,y)}{P(x)P(y)}$

Sampling distribution
 Variability of summaries of random samples from the population (mean, variance)
 1. y_1, y_2, \dots, y_n collect from population
 2. calculate mean \bar{y}
 3. repeat 1, 2 large times
 4. The histogram of large value of \bar{y} approximates the sampling distribution of \bar{Y} .

CLT central
 y_1, y_2, \dots, y_n follow arbitrary probability distribution with expected value μ and sd σ and n is large.
 then $\bar{Y} \sim N(\mu, \frac{\sigma}{\sqrt{n}})$

- Sampling distribution
 Variability of summaries of random samples from the population (mean, variance)
 1. y_1, y_2, \dots, y_n collect from population
 2. calculate mean \bar{y}
 3. repeat 1, 2 large times
 4. The histogram of large value of \bar{y} approximates the sampling distribution of \bar{Y} .

CLT central
 y_1, y_2, \dots, y_n follow arbitrary probability distribution with expected value μ and sd σ and n is large.
 then $\bar{Y} \sim N(\mu, \frac{\sigma}{\sqrt{n}})$

- difference between standard deviation and standard error.
 sd: quantifies the variability in the population
 SE: uncertainty in a parameter of population
 sd remains the same when sample size \uparrow
 while SE \downarrow
 $SE = \frac{sd}{\sqrt{n}}$

Linear regression (discriminative classifiers)

$Y_i \sim N(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}, \sigma^2)$ $Y = X\beta + \epsilon$ matrix form
 $Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i, i=1, \dots, n$ $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$
 - β_1 describes change in mean response per unit increase in x_1 when x_2 is held constant
 - β_2 describes change in mean response per unit increase in x_2 when x_1 is held constant

Interaction model
 $Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i1} x_{i2} + \epsilon_i, i=1, \dots, n$ $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$
 - change in x_1 when $x_2 = x_0$ $\Delta Y = \beta_1 + \beta_3 x_0$
 - change in x_2 when $x_1 = x_0$ $\Delta Y = \beta_2 + \beta_3 x_0$

Rate of change due to one variable affected by the other.

- Least squares Estimation same as maximum likelihood estimation
 - minimize $(Y - X\beta)^T (Y - X\beta)$
 $\hat{\beta} = (X^T X)^{-1} X^T Y$
 $\hat{Y} = X \hat{\beta} = X (X^T X)^{-1} X^T Y = H Y$
 - residuals
 $e = Y - \hat{Y} = (I - H) Y$
- Partition sums of squares
 - Total sums of squares
 $SSTO = \sum (Y_i - \bar{Y})^2 = \sum (Y_i - \hat{Y}_i + \hat{Y}_i - \bar{Y})^2 = \sum (\hat{Y}_i - \bar{Y})^2 + \sum (Y_i - \hat{Y}_i)^2$
 $= SSR + SSE$
 $R^2 = \frac{SSR}{SSTO} = 1 - \frac{SSE}{SSTO}$
 R^2 is not that useful.
- Mean squared Error (MSE)
 $MSE(\hat{\theta}) = E[(\hat{\theta} - \theta)^2]$
 $= E[(\hat{\theta} - E(\hat{\theta})) + (E(\hat{\theta}) - \theta)]^2$
 $= E[(\hat{\theta} - E(\hat{\theta}))^2] + E[(E(\hat{\theta}) - \theta)^2] + 2(E(\hat{\theta}) - \theta)E(\hat{\theta} - E(\hat{\theta}))]$
 $= Var(\hat{\theta}) + Bias^2$
 $\therefore E(\hat{\theta} - E(\hat{\theta})) = E(\hat{\theta}) - E(\hat{\theta}) = 0$
 $\therefore E(\hat{\theta}) - \theta$ is constant we drop the cross term
 $MSE(\hat{\theta}) = Var(\hat{\theta}) + Bias^2$

- Qualitative predictors
 $Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \epsilon_i$
 $x_{i1} = 1$ if stock firm
 $x_{i1} = 0$ if foreign firm
 x_{i2} and x_{i3} terms more flexible, interaction model.
- Three groups
 $Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \epsilon_i$
 $x_{i1} = 1$ if stock firm
 $x_{i1} = 0$ if foreign firm
 x_{i2} and x_{i3} terms more flexible, interaction model.

- Gauss-Markov theorem
 Least squares estimator has the smallest mean squared error of all linear estimators with no bias.
- Model selection and bias-variance tradeoff
 - Loss function $L(Y, \hat{Y})$
 - risk $R = E[L(Y, \hat{Y})] = \sum_{y,y'} L(y, y') P(y, y')$
 - estimated risk $R = \frac{1}{n} \sum_{i=1}^n L(Y_i, \hat{Y}_i)$
 $Y = f(X) + \epsilon, L = (Y - \hat{Y})^2$
 $E[(Y - \hat{Y})^2 | X = x] = \sigma^2 + [f(x) - \frac{1}{n} \sum_{i=1}^n f(x_i)]^2 + \frac{\sigma^2}{n}$
 $= \text{irreducible error} + \text{Bias}^2 + \text{Variance}$
- AIC and BIC criteria to choose between models
 $AIC = \frac{SSE}{n} + 2p$ or $\log \frac{SSE}{n} + 2p$
 $BIC = \frac{SSE}{n} + p \log(n)$ (Cheverson penalty) $\log \frac{SSE}{n} + p \log(n)$
- cross-validation
 Iteratively we each part for training/variable selection/validation
 cross-validation should be done as part of variable selection.

Ridge Regression
 - motivation: difficult inverting $X^T X$ when multicollinearity
 - minimize: $RSS + \lambda \sum \beta_j^2$ use CV to choose λ
 - $\hat{\beta} = (X^T X + \lambda I)^{-1} X^T Y$
 - biased but stable
 - scaled version of least squares $\beta / (C/H \lambda)$
 - elastic net penalty $\lambda \sum \beta_j^2 + \lambda \sum | \beta_j |$
 - $\lambda \rightarrow \infty$
 - training residual sum of squares \uparrow
 - test residual \downarrow then \uparrow
 - variance \downarrow
 - squared bias \uparrow

Logistic regression discriminative classifiers
 - $E(Y_i) = g(\beta_0 + \beta_1 x_i)$ or $g^{-1}(E(Y_i)) = \beta_0 + \beta_1 x_i$
 - g : mean response function g^{-1} : link function
 - $g(x)$: identity \rightarrow linear regression "does not have $0 \leq E(Y_i) \leq 1$ "
 - $\phi(x) \rightarrow$ probit \rightarrow non-normal distribution of E
 - $\frac{E(Y)}{H(Y)} \rightarrow$ logistic
 - motivation for probit regression: latent variable
 - $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$
 - $P(Y_i = 1) = P(\epsilon_i > -\beta_0 - \beta_1 x_i) = P(\epsilon_i < \beta_0 + \beta_1 x_i)$
 - sigmoidal response function
 - $E(Y_i) = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}$
 - logit link function odds
 - $\log\left(\frac{E(Y_i)}{1 - E(Y_i)}\right) = \log\left(\frac{P}{1 - P}\right) = \beta_0 + \beta_1 x_i$

confusion matrix
 predicted class
 True - True Neg False Pos N
 class + False Neg True Pos P
 \rightarrow $\frac{TP}{N}$ 1-specificity
 Type II error
 ROC curve, Receiver Operating Characteristic
 simultaneously displaying the two types of errors for all possible thresholds
 Y: true pos sensitivity X: false pos 1-specificity
 AUC = $2 \log L(D) + p$ BIC = $-2 \log L(D) + p \log(n)$
 - statistical regularization
 For logistic regression, maximize a penalized version
 Lasso $\min_{\beta} \left\{ \sum_{i=1}^n y_i (\beta_0 + \beta^T x_i) - \log \left(\sum_{i=1}^n e^{\beta_0 + \beta^T x_i} \right) - \lambda \sum_{j=1}^p |\beta_j| \right\}$

Nearest centroids
 (diagonal - covariance LDA)
 For $p \gg n$, covariance estimation is unstable
 assume independent features (diagonal Σ)
 - $\delta_k(x_i) = -\frac{1}{2} (x_i - \mu_k)^T \Sigma^{-1} (x_i - \mu_k) + \log \pi_k$
 - $\delta_k(x_i) = -\frac{1}{2} \frac{(x_{ij} - \bar{x}_{jk})^2}{s_j^2} + 2 \log \pi_k$
 for class k: $\bar{x}_{jk} = \sum_{i: i \in k} x_{ij} / N_k$
 $\bar{x}_j = \sum_{i=1}^N x_{ij} / N$
 observations $i = 1, \dots, N$
 features $j = 1, \dots, p$
 $\hat{Y}(x_i) = \arg \max_k \delta_k(x_i)$
 - The diagonal-covariance LDA classifier is equivalent to nearest centroid classifier after standardization, correcting for class prior probability

SVD singular value decomposition
 $X = U \Lambda V^T$
 NAP NAP PP PP
 $XV = U \Lambda$
 - scores
 The position of each observation in the new coordinate system of principal components
 $\frac{p}{N} (x_{ik} - \bar{x}_k) \cdot U_{kj} \rightarrow$ coordinate of observation i in direction j
 - Loadings
 weight U_{kj} of feature k in the direction j
 Scores and Loadings cannot be used as "measures of statistical significance".
 PCA of X span covariance matrix is PCA (p>n) when
 $\hat{\sigma}_{jk}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k)$
 $(N-1) \Sigma = X^T X = V \Lambda V^T = V \Lambda^* V^T$
 - eigenvalues of covariance matrix = $\frac{1}{N-1} \lambda_j^2$

Lasso Regression
 - minimize: $RSS + \lambda \sum |\beta_j|$
 - sign thresholding
 - weaker penalty on β than Ridge.
 - optimum is on the intersection with some axes

MLE independent but not identically distributed.
 $Y_i \sim \text{Bernoulli}(\pi_i)$ $\pi_i = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}$
 $f(Y_i) = \pi_i^{Y_i} (1 - \pi_i)^{1 - Y_i}$
 $\log L = \log \prod_{i=1}^n \pi_i^{Y_i} (1 - \pi_i)^{1 - Y_i}$
 $= \sum_{i=1}^n Y_i \log(\pi_i) + \sum_{i=1}^n (1 - Y_i) \log(1 - \pi_i)$
 $= \sum_{i=1}^n Y_i \log\left(\frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}\right) + \sum_{i=1}^n (1 - Y_i) \log\left(\frac{1}{1 + e^{\beta_0 + \beta_1 x_i}}\right)$
 $= \sum_{i=1}^n Y_i (\beta_0 + \beta_1 x_i) - \sum_{i=1}^n \log(1 + e^{\beta_0 + \beta_1 x_i})$
 $Y_i \sim \text{Binomial}(N_i, \pi_i)$
 $\log L = \log \prod_{i=1}^n \binom{N_i}{Y_i} \pi_i^{Y_i} (1 - \pi_i)^{N_i - Y_i}$
 $= \sum_{i=1}^n \left[Y_i \log(\pi_i) + (N_i - Y_i) \log(1 - \pi_i) + \log \binom{N_i}{Y_i} \right]$
 $= \sum_{i=1}^n \left[Y_i \log\left(\frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}\right) + (N_i - Y_i) \log\left(\frac{1}{1 + e^{\beta_0 + \beta_1 x_i}}\right) + \log \binom{N_i}{Y_i} \right]$
 $= \sum_{i=1}^n Y_i (\beta_0 + \beta_1 x_i) - \sum_{i=1}^n \log(1 + e^{\beta_0 + \beta_1 x_i}) + \sum_{i=1}^n \log \binom{N_i}{Y_i}$
 $= \log L_{\text{binomial}} + \text{constant}$
 $\log L$ binomial lead same parameter estimates and inferences but different deviances
 $D(y) = -2 (\log(P(y|\hat{\theta})) - \log(P(y|\hat{\theta}^*)))$

* Variable selection should be done as part of cross-validation
 $Y | X \sim \text{Multinomial}(N_i, \pi_1(x), \dots, \pi_K(x))$
 $\log \frac{\pi_j(x)}{\pi_i(x)} = \alpha_j - \alpha_i + \beta^T x$
 $\pi_j(x) = \frac{e^{\alpha_j + \beta^T x}}{1 + \sum_{k=2}^K e^{\alpha_k + \beta^T x}}$
 $\pi_j(x) = \frac{e^{\alpha_j + \beta^T x}}{1 + \sum_{k=2}^K e^{\alpha_k + \beta^T x}}$
 for $j=2, \dots, K$ \rightarrow ordinary logistic regression

Comments on PCA
 - Decomposition of Σ helps interpretation
 $\frac{1}{N-1} \sum_{i=1}^N \text{sample variance of columns of } X$
 $\frac{1}{N-1} \sum_{i=1}^N (x_{ij} - \bar{x}_j)^2 = \frac{1}{N-1} \text{tr}(X^T X)$
 $= \frac{1}{N-1} \text{tr}(V \Lambda^* V^T) = \frac{1}{N-1} \text{tr}(\Lambda^* V V^T)$
 $= \frac{1}{N-1} \text{tr}(\Lambda^*) = \sum_{j=1}^p \lambda_j^2$
 - The proportion of total variance explained by the principle component j is $\frac{\lambda_j^2}{\sum \lambda_k^2}$
 - doesn't distinguish useful and nuisance variation between-group within-group
 - can be driven by large nuisance variation
 - can not use scores as evidence of the ability of the features
 - PCA performs best when all features have similar nuisance variation

Inversion of Σ^{-1} in $\delta(x)$
 For LDA discriminant function is
 $\delta_k(x) = -\frac{1}{2} \log |\Sigma| - \frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log \pi_k$
 $\Sigma^{-1} = (V \Lambda^* V^T)^{-1}$
 $\therefore (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) = [V^T (x - \mu_k)]^T [\Lambda^* V^T (x - \mu_k)]$
 ① sphere the data $V^T (x - \mu_k) \rightarrow x^*$ $V^T \mu_k \rightarrow \mu^*$
 ② classify to closest class centroid.

Fisher's approach
 project high-dimensional data onto a lower-dimensional space to best separate class
 $\max \frac{\text{between-group variance}}{\text{within-group variance}} = \max_a \frac{a^T B a}{a^T S a}$
 - transform the system coordinates s.t. $a^T S a = 1$
 - maximize $a^T B a$ in the transformed system

Regularized LDA regularization
 $\hat{\Sigma}_k(a) = \hat{\Sigma}_k + (1 - \alpha) \hat{\Sigma}$ shrink the separate covariances of LDA toward a common covariance as in LDA
 $\hat{\Sigma}(a) = \alpha \hat{\Sigma} + (1 - \alpha) \hat{\sigma}^2 I$
 - shrink the common covariance in LDA toward a scalar covariance

MLE
 $l_i = \log \prod_{j=1}^J \pi_j(x_i) y_{ij}$
 $= \sum_{j=1}^J y_{ij} \log \pi_j(x_i) + (1 - y_{ij}) \log \pi_k(x_i)$
 $= \sum_{j=1}^J y_{ij} \log \frac{e^{\beta_j(x_i)}}{1 + \sum_{k=1}^K e^{\beta_k(x_i)}} + \log \pi_k(x_i)$
 $= \sum_{j=1}^J y_{ij} (\alpha_j + \beta_j x_i) - \log \left(1 + \sum_{k=1}^K e^{\alpha_k + \beta_k x_i} \right)$
 maximize $\sum_{i=1}^n l_i$ with respect to α_j and β_j
 can fit J-1 logistic regressions for J-1 response categories or baseline (the same model)

Generative vs discriminative
 $P(Y|X) = \frac{P(Y) \cdot P(X|Y)}{P(X)}$
 generative classifiers
 - specify prior probability $P(Y)$
 - assume conditional distribution $P(X|Y)$
 - use Bayes rule to derive the posterior $P(Y|X)$
 example LDA
 discriminative
 - estimate posterior $P(Y|X)$ directly
 example: linear regression logistic regression

connection of LDA and logistic regression
 - linear decision boundaries
 - logistic regression estimated using MLE
 - coefficients in LDA are estimated using $\hat{\mu}$ and $\hat{\sigma}$ from MVM
 $P(Y=k|X) = \frac{\delta_k(x) \pi_k}{\sum_{k=1}^K \delta_k(x) \pi_k}$
 MAP decision
 $\hat{Y}(x) = \arg \max_k P(Y=k|X) = \arg \max_k \delta_k(x) \pi_k$
 $\delta_k(x)$ types:
 - Gaussian, same $\Sigma \rightarrow$ LDA
 - Gaussian, different $\Sigma_k \rightarrow$ LDA
 - mixture of Gaussians \rightarrow Mixture models
 - $\prod_{j=1}^p \delta_k(x_j)$ any variable independent Naive Bayes
 - $\delta_k(x)$ nonparametric \rightarrow Kernel estimates.

Nearest regular centroid regularization
 $d_{ik} = \tilde{x}_{ik} - \tilde{x}_k$ for feature i class k
 $d_{ik} = \frac{\tilde{x}_{ik} - \tilde{x}_k}{m_k - 1}$
 $S_i^2 = \frac{1}{n-k} \sum_{k=1}^K \sum_{i \in k} (x_{ij} - \tilde{x}_{jk})^2$
 $m_k = \sqrt{1/n_k - 1/n}$
 - Nearest shrunken centroids
 start with nearest centroids
 $\tilde{x}_{ik} = \tilde{x}_i + m_k \cdot S_i \cdot d_{ik}$
 $\tilde{x}_{ik} = \tilde{x}_i + m_k \cdot S_i \cdot d_{ik}$ shrink d_{ik} to 0
 $d_{ik} = \text{sign}(d_{ik}) (|d_{ik}| - \Delta)_+$
 Δ is chosen by CV $\propto (\log\text{-likelihood or error rate})$

Soft thresholding vs hard
 $d_{ik}' = d_{ik} \cdot I(|d_{ik}| > \Delta)$
 we prefer soft-thresholding as it is a smoother operation and works better and not jumpy.
 discriminant score
 $\delta_k(x^*) = \sum_{i=1}^p \frac{(x_i^* - \tilde{x}_{ik})^2}{S_i^2} - 2 \log \pi_k$
 \uparrow
 standardized squared distance of x^* to the k th shrunken centroid.

Nearest centroids
 $\hat{Y}(x) = \arg \max_k P(Y=k|X) = \arg \max_k \delta_k(x) \pi_k$
 $\delta_k(x)$ types:
 - Gaussian, same $\Sigma \rightarrow$ LDA
 - Gaussian, different $\Sigma_k \rightarrow$ LDA
 - mixture of Gaussians \rightarrow Mixture models
 - $\prod_{j=1}^p \delta_k(x_j)$ any variable independent Naive Bayes
 - $\delta_k(x)$ nonparametric \rightarrow Kernel estimates.

LDA one predictor
 $\delta_k(x) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{1}{2\sigma^2} (x - \mu_k)^2\right)$
 $\hat{Y}(x) = \arg \max_k P(Y=k|X)$
 $= \arg \max_k \left[x \cdot \frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k) \right]$
 $= \arg \max_k \delta_k(x)$
 - MVM distribution $B \Sigma$
 $\hat{Y}(x) = \arg \max_k P(Y=k|X)$
 $= \arg \max_k \frac{\delta_k(x) \pi_k}{\sum_{k=1}^K \delta_k(x) \pi_k}$
 $= \arg \max_k \delta_k(x) \pi_k$
 $= \arg \max_k \left[\log(\delta_k(x)) + \log(\pi_k) \right]$
 $= \arg \max_k \left[-\log(2\pi) - \frac{1}{2\sigma^2} (x - \mu_k)^2 - \frac{1}{2\sigma^2} \mu_k^2 + \log(\pi_k) \right]$
 $= \arg \max_k \left[-\frac{1}{2\sigma^2} (x - \mu_k)^2 - \frac{1}{2\sigma^2} \mu_k^2 + \log(\pi_k) \right]$
 $= \arg \max_k \left[x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log(\pi_k) \right]$

decision boundary for LDA
 $P(Y=k|X) = P(Y=l|X)$
 $= x^T \log \frac{\pi_k}{\pi_l} - \frac{1}{2} (\mu_k - \mu_l)^T \Sigma^{-1} (\mu_k - \mu_l) + x^T \Sigma^{-1} (\mu_k - \mu_l) = 0$
 - QDA or LDA Quadratic Discriminant Analysis
 $\hat{Y}(x) = \arg \max_k \delta_k(x) \pi_k = \arg \max_k \left[-\log(2\pi) - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) - \frac{1}{2} \mu_k^T \Sigma_k^{-1} \mu_k + \log(\pi_k) \right]$
 $(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log(\pi_k)$
 - decision boundary
 $X: \log \frac{\pi_k}{\pi_l} - \log \frac{|\Sigma_l|}{|\Sigma_k|} - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) - \frac{1}{2} (x - \mu_l)^T \Sigma_l^{-1} (x - \mu_l) = 0$

Estimate the decision boundary
 $\hat{\mu}_k = \mu_k / N$
 $\hat{\mu}_k = \sum_{i: Y_i=k} X_i / N_k$
 $\hat{\Sigma} = \sum_{k=1}^K \sum_{i: Y_i=k} (X_i - \hat{\mu}_k)(X_i - \hat{\mu}_k)^T / (N - K)$
 - $(K-1)(p+1)$ parameters for LDA
 - $(K-1)(p(p+1)/2 + 1)$ parameters for QDA

likelihood vs probability
 - Probability is used before data are available to describe future outcomes given a fixed value for the parameter.
 - Likelihood is used after data are available to describe a function of a parameter for a given outcome.
 * Why not least squares in logistic regression?
 In logistic regression, the errors are not expected to have the same variance: should have high var when p near 0 or 1. Low variance towards extremes. Leads to IRWLs but not simple LS. to iteratively weighted least squares