Logistic regression vs LDA

LDA $\quad \log \dfrac{Pr[G=0|X=x]}{Pr[G=1|X=x]} = \log \dfrac{\widehat{\pi_0}}{\widehat{\pi_1}} - \dfrac{1}{2}(\mu_0 + \mu_1)^T \Sigma^{-1}(\mu_0 - \mu_1)$

$$+ x^T \Sigma^{-1}(\mu_0 - \mu_1)$$

$$= \alpha_{01} + \alpha_1^T x$$

similar to the logistic regression

$$= \beta_{01} + \beta_1^T x$$

$$Pr[G=k, X=x] = Pr(x) \; Pr[G=k|X=x]$$

both LDA and logistic regress.

LDA also model this part

$$Pr(x) = \sum_{k=1}^{K} \widehat{\pi_k} \, \phi(x; \mu_k, \Sigma)$$

Gaussian

$$\dfrac{e^{\beta_{0k} + \beta_k^T x}}{1 + \sum_{\ell=1}^{K} e^{\beta_{0\ell} + \beta_\ell^T x}}$$

# Model assessment and model selection

① ②

① evaluate the performance (in term of prediction) of a selected model

② select the best models (for prediction)

GOAL of a prediction model

generalization : a prediction model must be valid in braod generality, and not valid for a specific dataset

# Bias, variance, model complexity

Define:

$Y$ = target variable

$X$ = input matrix

$\hat{f}(x)$ prediction rule, that is trained on a training set $\mathcal{T}$

Error is measured through a loss function

$$L\left(Y, \hat{f}(x)\right)$$

that should penalize differences between $Y$ and $\hat{f}(x)$

typical choice for continuous outcomes

$$L(Y, \hat{f}(x)) = \begin{cases} (Y - \hat{f}(x))^2 & \text{quadratic loss} \\ |Y - \hat{f}(x)| & \text{absolute loss} \end{cases}$$

Test error (generalization error) is a prediction error computed on an __independent__ sample

$$Err_{\mathcal{T}} = E\left[L(Y, \hat{f}(x)) \mid \mathcal{T}\right]$$

random, $X, Y$

$\mathcal{T}$ is fixed, is the specific training set on which we derive our prediction rule

In general, we would like to minimize the expected prediction error

$$Err = E\left[L(Y, \hat{f}(x))\right] = E\left[Err_{\mathcal{T}}\right]$$

We do not want a model with the smallest prediction error for the specific training set, but for the general case

Since we love our training set, we are going to estimate $Err_{\tau}$
  $\hookrightarrow$ the training error is $\underline{NOT}$ a good estimate of $Err_{\tau}$
      $\hookrightarrow \overline{err} = \frac{1}{N} \sum_{i=1}^{N} L(y_i, f(x_i))$

We do $\underline{NOT}$ want to minimize the training error.
We saw that by <mark>increasing the model complexity, we can always decrease the training error</mark>

OVERFITTING $\longrightarrow$ our model is specific for the training set
             $\searrow$ generalizes poorly

---

Similar story for categorical outcome                    <span style="color:red">indicator function</span>
$G$ : target variable $\longrightarrow$ takes $\underline{K}$ values in $\mathcal{G}$          <span style="color:red">$\mathbb{1}(x \neq x_0) \begin{cases} 1 & x \neq x_0 \\ 0 & x = x_0 \end{cases}$</span>

typical loss functions in this case    $L(G, \hat{G}(x)) \cdot \begin{cases} \underline{\mathbb{1}(G \neq \hat{G}(x))} & \text{0.1 loss} \\ \underbrace{-2 \log L(\beta)}_{\ell(\beta)} & \text{deviance} \end{cases}$

- $-2\,\ell(\beta)$ is a general loss-function, it can be used in all cases
  (binomial, Gamma, Poisson, log-normal, ...)

- the factor $-2$ is added to make the loss function be equal to the squared loss in the Gaussian case

$$L(\beta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma}\sum_{i=1}^{N}(y_i - x_i^T\beta)^2\right\}$$

$$\ell(\beta) = -\frac{1}{2}\underbrace{\sum_{i=1}^{N}(y_i - x_i^T\beta)^2}_{\color{red}\text{sum of squares}}$$

<span style="color:red">$-2\,\ell(\beta) = $ squared loss</span>

In an ideal situation
⤷ a lot of data
we can randomly split the observations in three
independent sets

| training | validation | test |
|----------|------------|------|

model selection        model assessment

- training set : contains the data on which we fit our models
- validation set : data we use to identify the best model
- test set : data to assess the performance of the selected model

⤷ this set must be considered only at the end of the analysis
i.e. only when we have already chosen the best model
Must be ignored for model selection
– avoid overoptimism

How to split the data in the three sets :
- no general rule
- book suggestion is :   50%   training set
                          25%   validation set
                          25%   test set

- it depends:
  + sample size
  + on the signal-to-noise ratio
  + complexity of the model we are considering

# The Bias-Variance Decomposition

$$Y = f(x) + \varepsilon \quad , \quad E[\varepsilon] = 0$$
$$Var[\varepsilon] = \sigma_\varepsilon^2$$

$$Err(x_0) = E\left[(Y - \hat{f}(x_0))^2 \mid X = x_0\right]$$

$$= \sigma_\varepsilon^2 + \left(E[\hat{f}(x)] - f(x_0)\right)^2 + E\left[(\hat{f}(x_0) - E[\hat{f}(x_0)])^2\right]$$

$$= \underset{error}{irreducible} + bias^2 + variance$$

where

- $\sigma_\varepsilon^2$, the variance of the target around the true mean, so we cannot do anything about that

- $bias^2$, the squared difference between the average of our estimates and the true mean

- variance, the expected squared difference between $\hat{f}(x)$ and its mean

## KNN

$$\hat{f}(x_0) = \frac{1}{K} \sum_{x_i \in N_k(x_0)} y_i$$

$$Err(x_0) = \sigma_\varepsilon^2 + \left[f(x_0) - \frac{1}{K}\sum_{\ell=1}^{K} f(x_{(\ell)})\right]^2 + \frac{\sigma_\varepsilon^2}{K}$$

$$E[\hat{f}(x)] = E\left[\frac{1}{K}\sum_{\ell=1}^{K} Y_{(\ell)}\right] = \frac{1}{K}\sum_{\ell=1}^{K} E[Y_{(\ell)}] = \frac{1}{K}\sum_{\ell=1}^{K} f(x_{(\ell)})$$

$$Var[\hat{f}(x_0)] = Var\left[\frac{1}{K}\sum_{\ell=1}^{K} Y_{(\ell)}\right] = \frac{1}{K^2}\sum_{\ell=1}^{K} Var[Y_{(\ell)}] = \frac{1}{K^2}\sum_{\ell=1}^{K} \sigma_\varepsilon^2 = \frac{K}{K^2}\sigma_\varepsilon^2 = \frac{\sigma_\varepsilon^2}{K}$$

$$\# K \propto \frac{1}{complexity}$$

increase K  ⟶ decrease complexity
           ⟶ more bias
           ⟶ reduce the variance

Similar for linear model $\hat{f}(x;\beta) = x^T\beta$
- slightly more difficult to derive, it turns out to be

$$\frac{1}{N}\sum_{i=1}^{N} Err(x_i) = \sigma_\varepsilon^2 + \frac{1}{N}\sum_{i=1}^{N}\left[f(x_i) - E[f(x_i)]\right]^2 + \frac{P}{N}\sigma_\varepsilon^2$$

$P\uparrow$

called the in-sample error          $P = \#$ of variables

increasing the model complexity,
we increase the variance
component of the error

For regularized regression (e.g., lasso, ridge,...) the form is the same,
but there is an additional dependence on the tuning (complexity) parameter
$\alpha$
• we can go more into details on the bias component

$$E_{x_o}\left[f(x_o) - E[\hat{f}(x_o)]\right]^2 = E_{x_o}\left[f(x_o) - x^T\hat{\beta}_*\right] + E_{x_o}\left[x_o^T\beta - E[x_o^T\hat{\beta}_\alpha]\right]^2$$

$$\hat{\beta}_* = \arg\min_{\beta} E\left[f(x) - x^T\beta\right]^2$$

Average [model bias]$^2$ + Ave [Estimation bias]$^2$

difference between the
true function and the best-fitting
linear approximation

error between the average
estimate and the best model

additional bias that we add
when, e.g., add shrinkage
$OLS = 0$

difference between truth and best model
- we can reduce it only increasing the model
space (more variables, more complex relationship - interaction -,...)