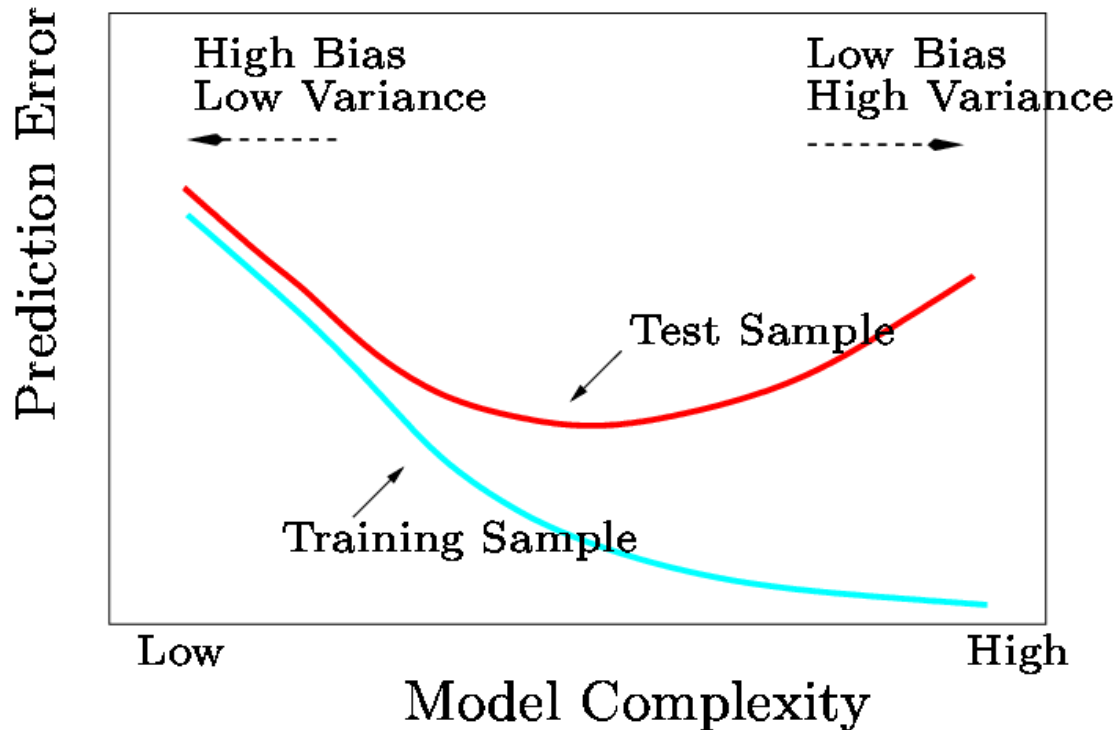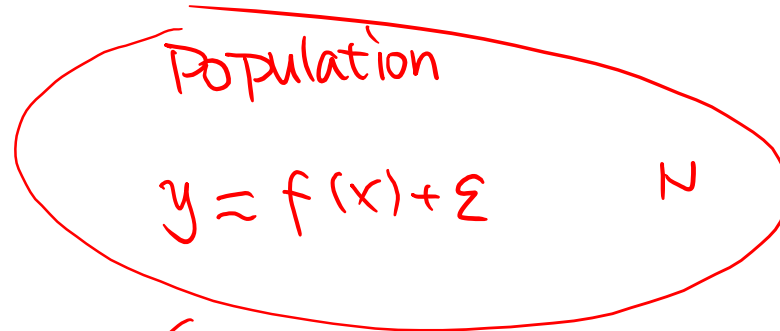# DSC5103 Statistics

Session 2. The K-Nearest Neighbor Algorithm

# Last time

- Out-of-sample prediction performance as the correct measure
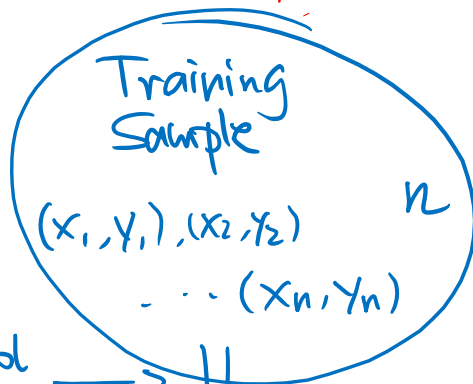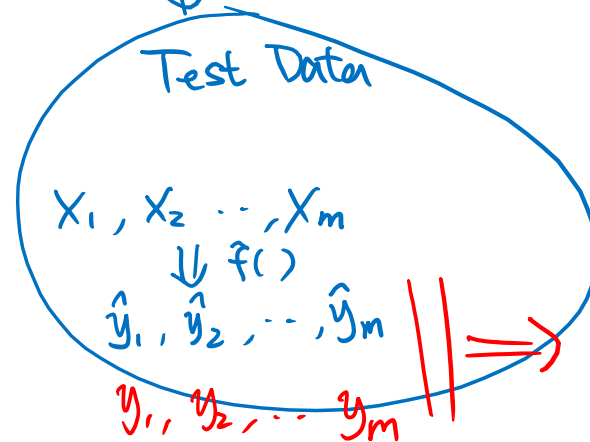- The Bias-Variance decomposition and trade-off

Population

$$y = f(x) + \varepsilon \qquad N$$

$\bigvee$

$\bigcup \longrightarrow$ Monte Carlo Simulation

Training Sample $\qquad n$

$(x_1, y_1), (x_2, y_2)$

$\cdots (x_n, y_n)$

tool $\longrightarrow$

$\hat{f}(\ )$

$\Downarrow$

Training Error

Test Data

$x_1, x_2 \cdots, x_m$

$\Downarrow \hat{f}(\ )$

$\hat{y}_1, \hat{y}_2, \cdots, \hat{y}_m \| \Rightarrow$ Test Error

$y_1, y_2, \cdots y_m \|$
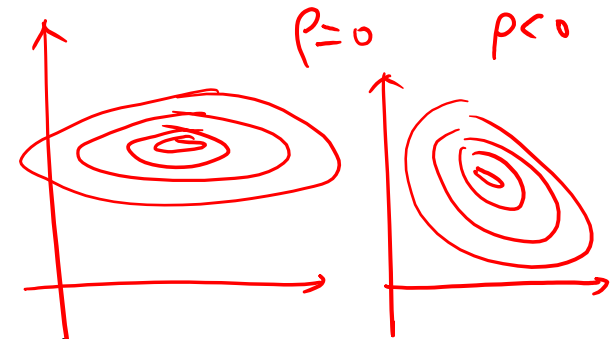
3

# Simulation Fundamentals

$f(\ )$

- Play the God's role and generate data with a known mechanism

- Stochastic models
  - Random variables, their distributions and parameters
  - The relationship among variables, dependency structures

- Monte Carlo simulation
  - Use computer to simulate outcomes of a stochastic model

  - To mimic the process of obtaining a sample from the population
  - By comparing the known population parameters and the estimates made from the sample, we can better evaluate our estimation methods

# Simulation in R

$\rho = 0$      $\rho < 0$

- Simulating random variables with a known distribution
  - Random sampling   *sample( )*
  - Uniform   [min, max]
  - Poisson   $\lambda$ (mean), $\{0, 1, 2, \cdots\}$
  - Binomial   $n, P$    $\{0, 1, 2, \cdots, n\}$
  - Normal   $u, \sigma^2$    $-\infty, +\infty$   $f(x)$    $\varepsilon \sim N(0, \sigma^2)$
  - Multivariate Normal   $u_W$     $W$    $\rho > 0$    $x$

$(u_H, u_W)$   $\begin{bmatrix} \sigma_H^2 & \rho\sigma_{H\sigma_W} \\ \rho\sigma_{H\sigma_W} & \sigma_W^2 \end{bmatrix}$    $u_H$   $H$

- Simulating stochastic models
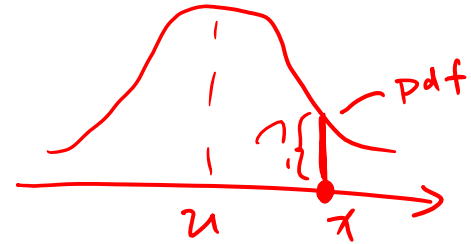  - A 1-D linear model   $y = \beta_0 + \beta_1 X + \varepsilon$   $\boxed{\beta_0}, \boxed{\beta_1}$   $\varepsilon \sim \underline{N}(0, \sigma^2), \boxed{\sigma^2}$
    $X \sim \underline{Unif}(a, b), \boxed{a}, \boxed{b}$
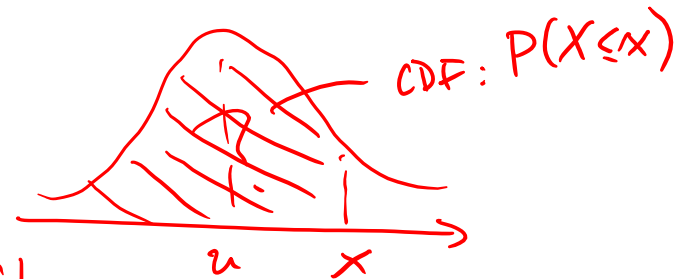  - A 2-D classification example

5

Prob. Functions in R: unif, binom, norm, pois, mvnorm

$\underset{=}{r}\underline{norm}(n, \overset{mean}{u}, \overset{sd}{\sigma})$: Simulate $n$ points
from the dist

$\underset{=}{d}norm(x, u, \sigma)$: PDF (Prob. density func)


pdf

$\cancel{p}norm(x, u, \sigma^2)$: CDF


CDF: $P(X \leq x)$

$q norm(p, u, \sigma^2)$: quantile

$$y: \{M, F\}$$
$$0, 1 \quad -1 \leq \rho \leq 1$$

$X_2 \; 70$

$(W) \; 60$

$160 \quad X_1 \; (H) \; 170$

$M: \quad (u_1^M, u_2^M), \quad \begin{bmatrix} \sigma_{1m}^2 & \rho \sigma_{1m} \sigma_{2m} \\ \rho \sigma_{1m} \sigma_{2m} & \sigma_{2m}^2 \end{bmatrix}$

$F: \quad (u_1^F, u_2^F), \quad \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix}$

$X_2$

$X_1$

$\times: M$

$0: F$

7

# K-Nearest Neighbors (KNN)

*[handwritten top-left: box with X, W, H, 10kg, 20cm]*

- <u>k Nearest Neighbors</u> is a flexible *nonparametric* approach for both regression and classification.

- For any given *X*, we find the k closest neighbors to *X* in the training data, and examine their corresponding *Y*. We use
  - the **average** of the neighbors' *Y* as prediction for regression;
  - the **majority votes** of the neighbors' *Y* as prediction for classification.

  *[handwritten: ⇒ Proportion ⇒ Prob]*

- The smaller that k is, the more flexible the method will be.

  *[handwritten: (High Var, Low Bias) k=1 |————— most flexible ——— k* ——— | k=n ($\hat{y} = \bar{y}$) (high Bias, low Var) Least flexible]*

- !!!BE CAREFUL!!!
  - How to define "closest"? How to measure distance in the space of *X*? *[handwritten: ⇒ Normalize]*
  - Categorical X?
  - Dimensionality??
  - Variable selection??? *[handwritten: ⇒ PCA]*
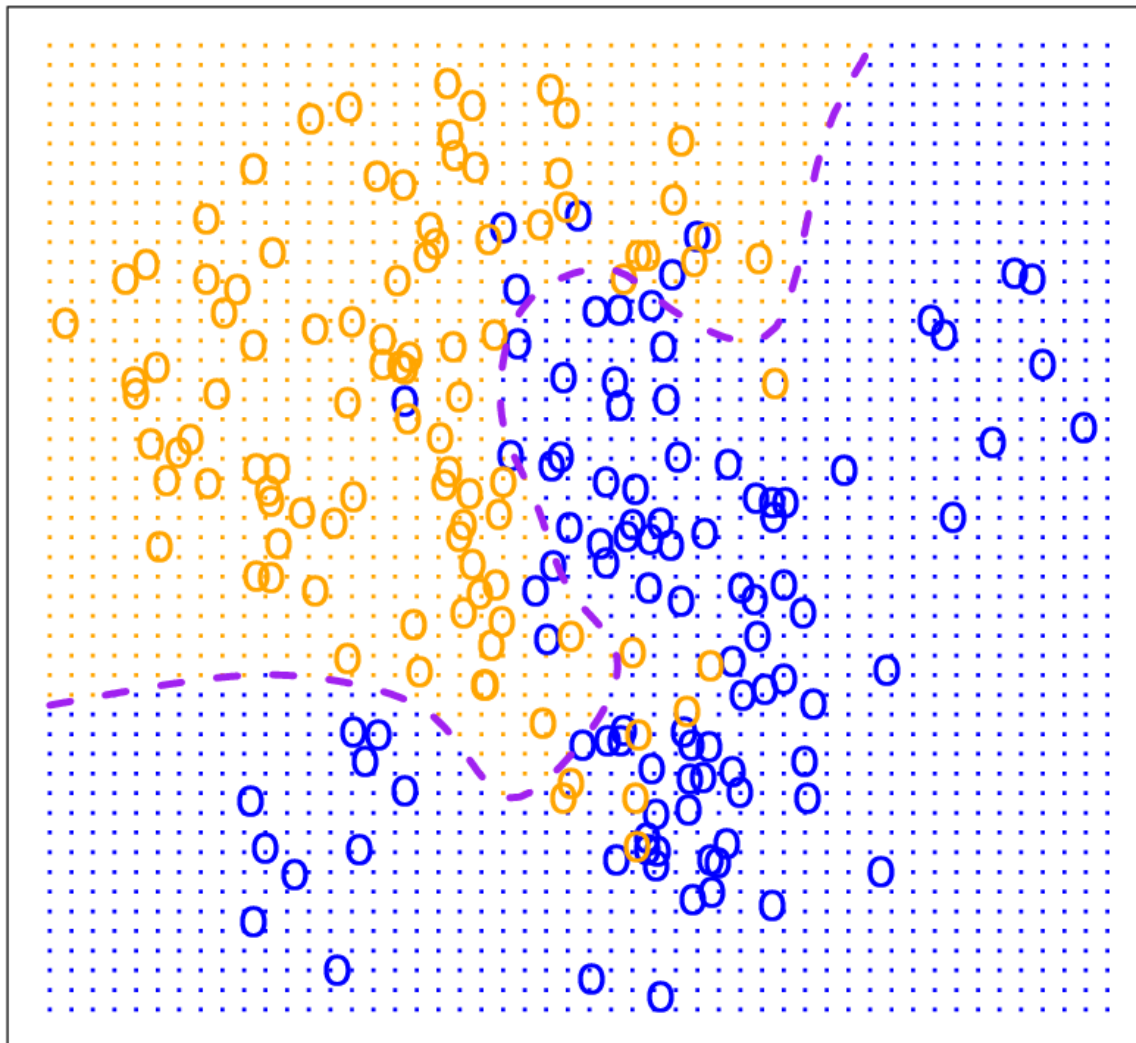
  *[handwritten: + Capture Nonlinearity]*

# KNN Example with k = 3

# Optimal Classifier

# Simulated Data: K = 10



KNN: K=10

# K = 1 and K = 100

# In-Sample vs. Out-of-Sample

# KNN Demo in R

- A regression problem using KNN


- Applying KNN on the Mixture Example
  - https://web.stanford.edu/~hastie/ElemStatLearn/datasets/mixture.example.info.txt

$y: \{0, 1\}$

$(99)$

$x_2$

$Px_1 \quad Px_2$

Grid: $69 \times 99$

$6831$

$x_1 \quad (69)$

- Homework 1
  - Learn RMarkdown (http://rmarkdown.rstudio.com/)

  - Test the curse of dimensionality of the KNN algorithm

- Next time:
  - Linear Regression