

Machine Learning I: Supervised Methods

B. Keith Jenkins

Announcements

- Slido event code: 53959875
- Homework 8 will be posted today after lecture (due Mon. 4/22)
- Project Report is due Friday 4/26
 - Templates have been posted
- Reminder: Final exam is Wed. 5/1/2024
 - 4:30 - 6:30 PM PDT
 - We will have a larger room! (SGM)

Today's lecture

- Density estimation techniques for machine learning (non-parametric)

- Preliminaries
- Convergence
- 2 approaches (KDE, kNN)
- Kernel density estimation (KDE)
- k-nearest neighbors density estimation (kNN)
- Discriminative and Generative approaches to ML
- Classification based on density estimation
- Generative (using KDE or kNN)
- Discriminative (using kNN or binary-window KDE)

DENSITY ESTIMATION TECHNIQUES FOR MACHINE LEARNING

Suppose we want to use a Bayes min. error classifier:

$$g_i(\underline{x}) = \ln \left[p(\underline{x} | S_i) P(S_i) \right], \quad i=1, 2, \dots, C$$

$$\text{or } \tilde{g}_i(\underline{x}) = \ln P(S_i | \underline{x})$$

but we don't know $p(\underline{x} | S_i)$, $P(S_i | \underline{x})$, maybe $P(S_i)$.

⇒ Estimate $p(\underline{x} | S_i)$ (or $P(S_i | \underline{x})$) from the data.

Or in regression problems, we will typically want to know $p(\underline{x}, y)$, $p(y | \underline{x})$, and/or $p(\underline{x})$. If they are not known:

⇒ Estimate $p(\underline{x}, y)$, $p(y | \underline{x})$, and/or $p(\underline{x})$ from the data.

Intro. & Assumptions

Suppose we want to est. $p(\underline{x} | S_1)$, $p(\underline{x} | S_2)$, ..., $p(\underline{x} | S_c)$

Assume data pts. of S_j are not useful in estimating $p(\underline{x} | S_i)$, $i \neq j$.

⇒ Estimate $p(\underline{x} | S_i)$ using only data pts. labeled S_i ,
for $i=1, 2, \dots, C$.

Notation: when considering $p(\underline{x} | S_i)$ only, drop $| S_i$ and class subscripts & superscripts for now.

\Rightarrow Estimate $p(\underline{x})$, given a set of data pts. (drawn i.i.d. from $p(\underline{x})$): $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N$.

Or, e.g., estimate $p(\underline{x}, y)$ from the data: (\underline{x}_i, y_i) , $i=1, \dots, N$.

GENERAL APPROACH TO ESTIMATE $p(\underline{x})$ FROM DATA

Vector \underline{x} : $P\{\underline{x} \text{ lies in } R\} = P = \int_R p(\underline{x}') d\underline{x}'$

Draw N vectors $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N$
i.i.d. from $p(\underline{x})$:

$$P_k = P\{k \text{ vectors lie in } R\} = ?$$

$$P_k = \binom{N}{k} P^k (1-P)^{N-k} \quad \text{Binomial}$$

$$\text{mean: } E\{k\} = NP = \bar{k} \Rightarrow P = \frac{\bar{k}}{N}$$

\Rightarrow reasonable est. for P is $\hat{P} = \frac{k}{N}$



$$\hat{P} = ? = \frac{4}{15}$$

Assume R is small, has volume V , has \underline{x} at its center.
Also assume $p(\underline{x})$ is continuous at \underline{x} .



$$P = \int_R p(\underline{x}') d\underline{x}' \approx p(\underline{x}) \int_R d\underline{x}' = p(\underline{x}) V$$

↑
approx. $p(\underline{x}') \approx p(\underline{x}) = \text{const. in } R.$

$$\Rightarrow p(\underline{x}) \approx \frac{P}{V} = \frac{\int_R p(\underline{x}') d\underline{x}'}{\int_R d\underline{x}'}$$

$$\text{use } P \approx \hat{P} = \frac{k}{N}$$

$$\hat{p}(\underline{x}) = \frac{\hat{P}}{V} = \frac{\frac{k}{N}}{V}$$

Annotations:

- $\#$ (training) data pts. that lie in R
- Total $\#$ of (training) data pts.
- Volume of R .

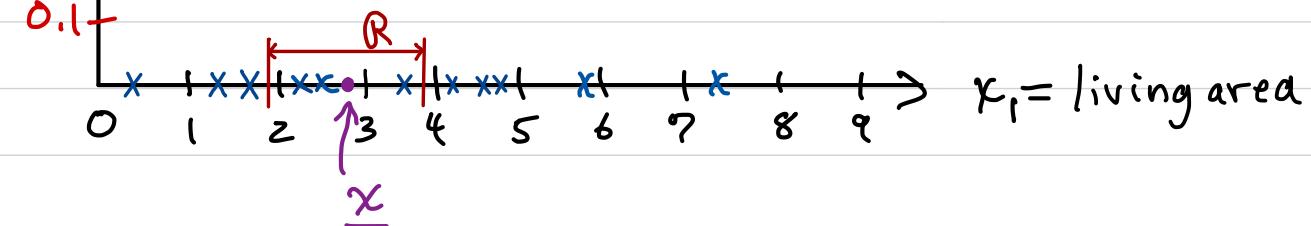
Let $\hat{p}(\underline{x}) \rightarrow \hat{p}(x_1 | S_1)$:

$$\hat{p}(x_1 | S_1) \quad \text{price increase}$$

\times data point

$$k=3, N=11, V=2$$

$$\Rightarrow \hat{p}(x_1 | S_1) = \frac{\frac{3}{11}}{2} = \frac{3}{22}$$



For theory, can construct a sequence

$$p_n(\underline{x}), \quad n=1, 2, \dots$$



Est. of $p(\underline{x})$ based on n data points.

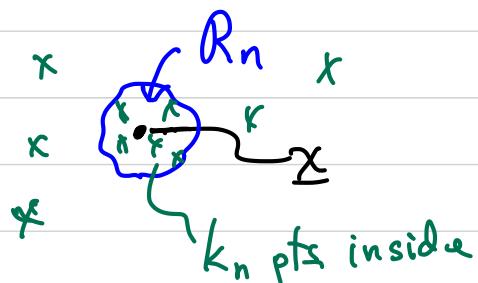
For each est. $p_n(\underline{x})$:

n data pts.

Region R_n containing \underline{x} , volume V_n

$k_n = \# \text{ pts. in } R_n$

$$\Rightarrow p_n(\underline{x}) = \frac{k_n}{n} / V_n$$



For the sequence $p_n(\underline{x}), n=1, 2, \dots$, to converge to $p(\underline{x})$ as $n \rightarrow \infty$, we need to satisfy:

$$1. \lim_{n \rightarrow \infty} V_n = 0$$

(spatial resolution $\rightarrow \infty$)

$$2. \lim_{n \rightarrow \infty} k_n = \infty \quad (\text{if } p(\underline{x}) \neq 0)$$

(quantization error $\rightarrow 0$)

$$\text{and } 3. \lim_{n \rightarrow \infty} \frac{k_n}{n} = 0$$

(% of data pts.
in $R_n \rightarrow 0$)

2 PRIMARY APPROACHES TO ESTIMATING $p(x)$

1. Specify $V_n = f(n)$

(perhaps also specify shape or profile of R_n)

$$\text{e.g., } V_n = \frac{1}{\sqrt{n}}$$

\Rightarrow Kernel density estimation (KDE) (also: Parzen Windows)
(PW)

2. Specify $k_n = f(n)$

$$\text{e.g., } k_n = \sqrt{n}$$

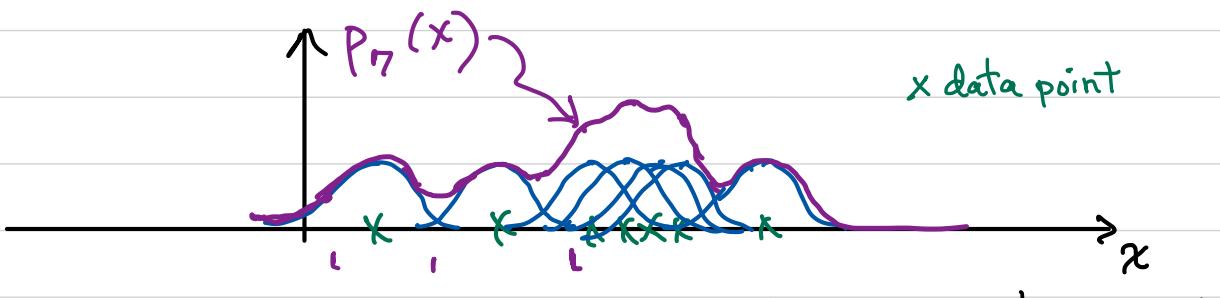
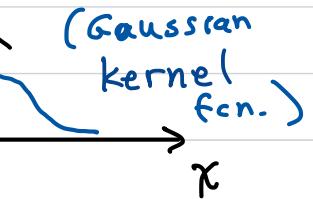
\Rightarrow k -Nearest Neighbors (k -NN)

Both can converge to $p(x)$ as $n \rightarrow \infty$ if V_n or k_n is chosen appropriately.

Kernel Density Estimation (KDE)

Ex:

- Create a window or kernel function
- Center the kernel fcn. at each data pt.



- Take sum of all kernel (window) fns.

Window fcn: $\Delta(\underline{u})$

$\Delta(\underline{x} - \underline{x}_i)$ centered at \underline{x}_i .

$$p_n(\underline{x}) = \frac{1}{n} \sum_{i=1}^n \Delta(\underline{x} - \underline{x}_i)$$

← KDE estimate

For $p_n(\underline{x})$ to be a density, require:

$$(*) \begin{cases} \int \Delta(\underline{u}) d\underline{u} = 1 \\ \Delta(\underline{u}) \geq 0 \end{cases}$$

If $\Phi(\underline{x})$ is an unnormalized window fcn.,

then we can define:

$$\Delta_n(\underline{x}) = \frac{1}{V_n} \Phi\left(\frac{\underline{x}}{h_n}\right)$$

Choose V_n to ensure $\int \Delta_n(\underline{x}) d\underline{x} = 1$. width parameter

Convergence of PW/KDE

Is $p_n(\underline{x})$ random or deterministic?

> $p_n(\underline{x})$ is random, because \underline{x}_i are drawn at random.

Does sequence of $p_n(\underline{x})$ converge to $p(\underline{x})$ as $n \rightarrow \infty$?

Convergence in mean square:

$$\text{Def: } \lim_{n \rightarrow \infty} \overline{p_n}(\underline{x}) = p(\underline{x})$$

$$\text{and: } \lim_{n \rightarrow \infty} \sigma_n^2(\underline{x}) = 0$$

$\sigma_n^2(\underline{x})$ is the variance of $p_n(\underline{x})$.

$$\text{If } \Delta_n(x) = \frac{1}{V_n} \mathbb{E}\left(\frac{x}{h_n}\right)$$

then the following conditions ensure convergence:

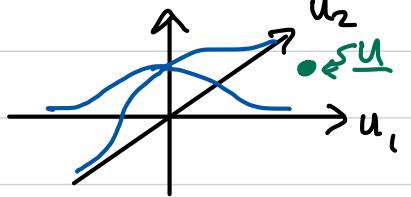
In addition to (*):

$$(1) \sup \Phi(u) < \infty$$

$$(2) \lim_{\|u\| \rightarrow \infty} \left[\Phi(u) \prod_{d=1}^D u_d \right] = 0$$

$\Phi(u)$ must be well behaved.

$\Phi(u)$



KDE convergence

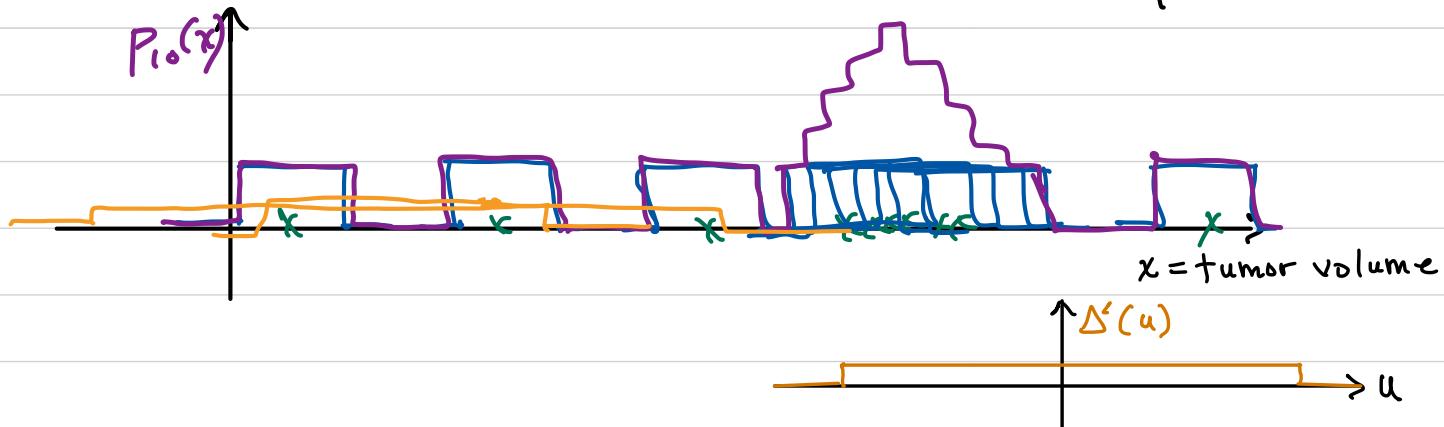
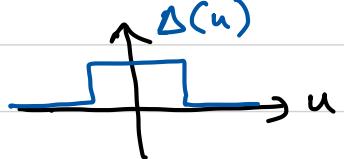
$$(3) \lim_{n \rightarrow \infty} V_n = 0$$

$V_n \rightarrow 0$, but not too fast.

$$\text{and } (4) \lim_{n \rightarrow \infty} n V_n = \infty$$

$$? V_n = \frac{1}{\sqrt{n}} \quad n V_n = \frac{n}{\sqrt{n}} = \sqrt{n}$$

(S) Ex: KDE with rectangular window fcn.



Ideally, would like window width to vary, dependent on local density of data points.

⇒ k-NN

k-NEAREST NEIGHBORS (k-NN) FOR DENSITY ESTIMATION

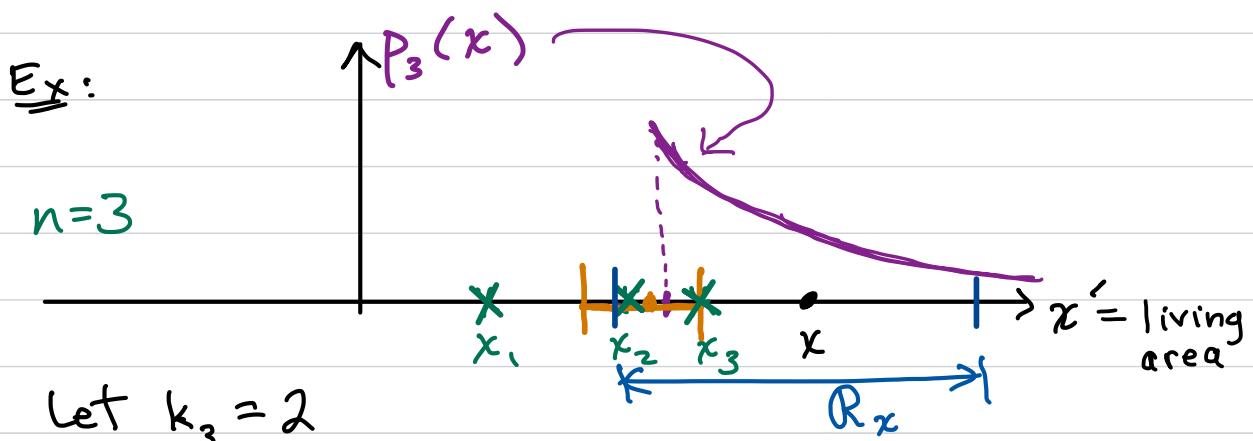
No width parameter (h_n).

No choice of window or kernel shape (in k-NN standard form)

→ Let region width vary according to local density of data points.

Center R_n at \underline{x} . Make R_n just large enough to capture k_n data pts. Specify $k_n = f(n)$.

Still: $P_n(\underline{x}) = \frac{\frac{k_n}{n} \leftarrow \text{specify const.}}{V_n \leftarrow \text{calculated based on } k_n}$



$$P_3(x) = \frac{\frac{k_3}{3}}{V_3} = \frac{\frac{2}{3}}{V_3(x)}$$

$$V_3(x) = 2|x - x_2|$$

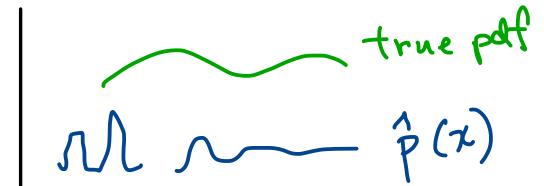
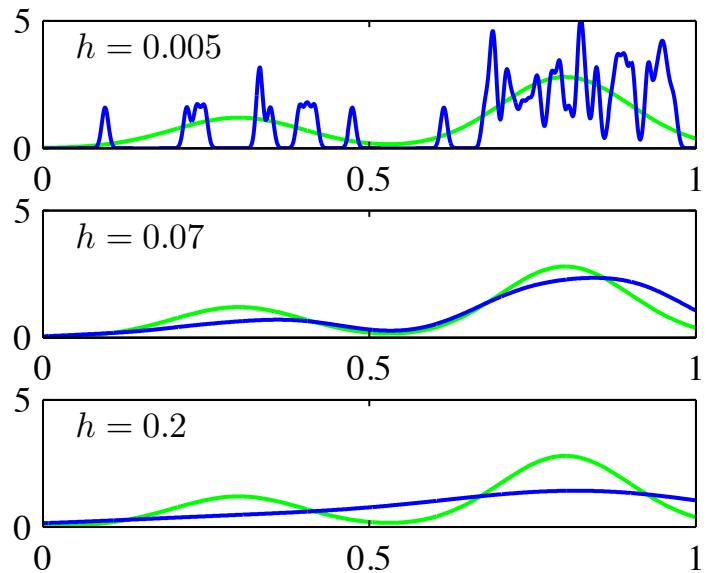
$$P_3(x) = \frac{\frac{2}{3}}{2|x - x_2|} = \frac{1}{3(x - x_2)}, \quad x \geq \frac{x_2 + x_3}{2}$$

If we choose $k_n = k_1 \sqrt{n}$, $k_1 > 0$, k_1 is a parameter
 Then $P_n(\underline{x})$ converges to $p(\underline{x})$ as $n \rightarrow \infty$.

Examples

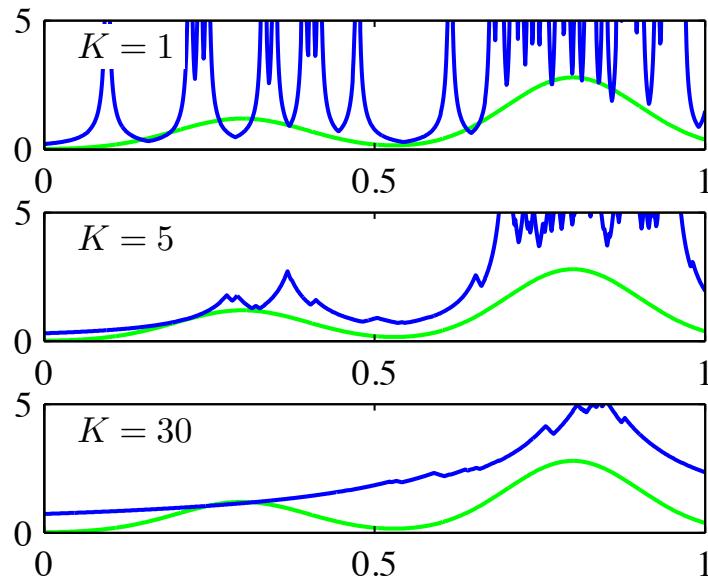
Dataset : 50 points drawn from pdf of green curve.

KDE



Bishop Fig. 2.25

kNN



Bishop Fig. 2.26

Generative and Discriminative Approaches for ML

1. Generative

→ Estimate $p(\underline{x}, y)$. (or, $p(\underline{x}|y)$ and $p(y)$). 

$$\text{From that, can get: } P(y|\underline{x}) = \frac{p(\underline{x}, y)}{p(\underline{x})}$$

Note: from $p(\underline{x}, y)$ we can synthetically generate more data.

2. Discriminative

→ Estimate $p(y|\underline{x})$ directly.

(Can't generate more data synthetically, unless also estimate $p(\underline{x})$.)

CLASSIFICATION BASED ON DENSITY ESTIMATION

> Option 1: Generative Approach

Apply above techniques to each class separately:

1. Estimate $p(\underline{x} | S_k)$ from data pts. $\underline{x}_i^{(k)}$ of S_k , separately for each class.

2. Use or estimate priors $P(S_k)$, $k=1, 2, \dots, C$

(S)

$$\text{e.g., } \hat{P}(S_k) = \frac{n^{(k)}}{N} \quad (\text{frequency of occurrence})$$

$n^{(k)} = \# \text{data pts. labeled } S_k \text{ (in training set)}$

$N = \# \text{data pts. in training set.}$

3. Use Bayes classifier

(min.error or min. risk)

$$\hat{p}(\underline{x} | S_k) \hat{P}(S_k) > \hat{p}(\underline{x} | S_j) \hat{P}(S_j) \quad \forall j \neq k$$

$$\Rightarrow \underline{x} \in S_k.$$

> Option 2 : Discriminative Approach

For the case of k-NN or KDE estimates using a binary-valued window fcn., we can derive an expression in advance based on estimates of posterior probabilities:

$$P(S_k | \underline{x}) > P(S_j | \underline{x}) \quad \forall j \neq k \Rightarrow \underline{x} \in S_k$$

as follows:

$$P_n(\underline{x}, S_i) = \frac{\frac{k_n^{(i)}}{n}}{V_n}$$

$k_n^{(i)}$ = # datapts. in R_n that belong to S_i .

n = total # datapts. over all classes

k_n = total # datapts. in R_n . ←

V_n = volume of R_n .

*{ k-NN est.
for classifier
specifies
this k. }*

$$P_n(\underline{x}, S_i) = P_n(S_i | \underline{x}) p_n(\underline{x})$$

$$P_n(\underline{x}) = \sum_{i=1}^C P_n(\underline{x}, S_i)$$

$$P_n(S_i | \underline{x}) = \frac{P_n(\underline{x}, S_i)}{\sum_{i'=1}^C P_n(\underline{x}, S_{i'})} = \frac{\frac{k_n^{(i)}}{n}}{\sum_{i'=1}^C \left(\frac{k_n^{(i')}}{n} \right)}$$

$$P_n(S_i | \underline{x}) = \frac{\sum_{i'=1}^C k_n^{(i')}}{\sum_{i'=1}^C k_n^{(i')}} = \frac{k_n^{(i)}}{k_n}$$

datapts in S_i that lie in R_n

datapts. over all classes that lie R_n .

$\Rightarrow P_n(S_i | \underline{x}) = \%$ of data pts. in R_n that belong to S_i .

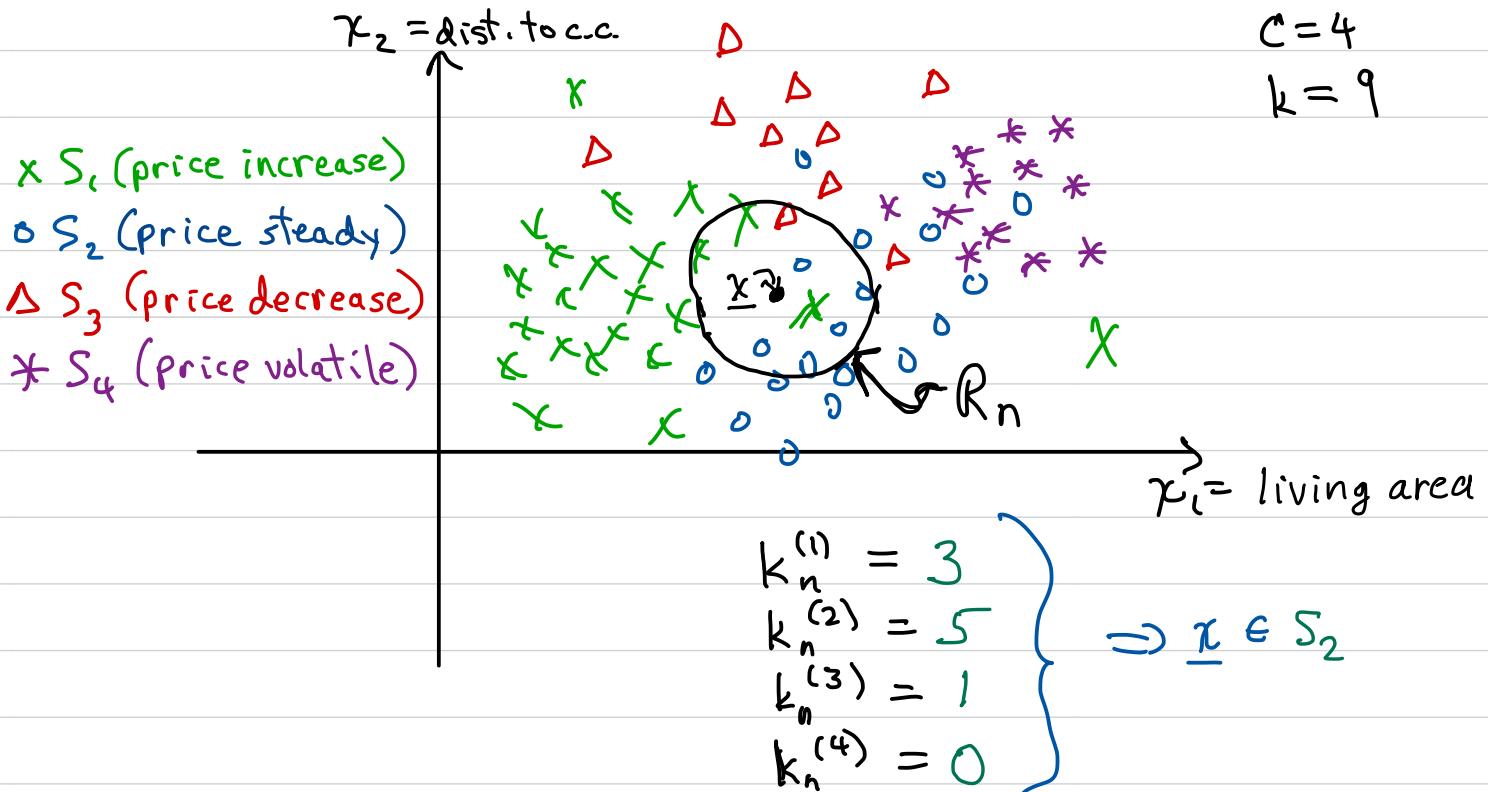
Use in a Bayes min.error classifier:

$$P(S_k | \underline{x}) > P(S_j | \underline{x}) \quad \forall j \neq k \Rightarrow \underline{x} \in S_k.$$

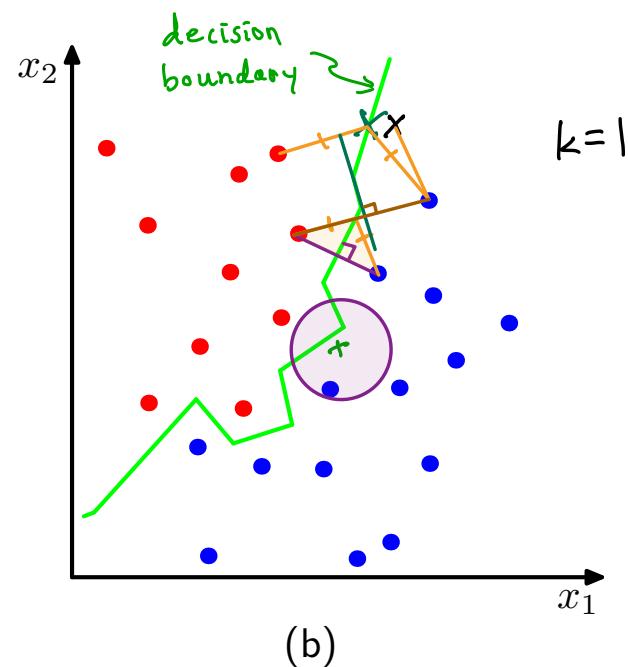
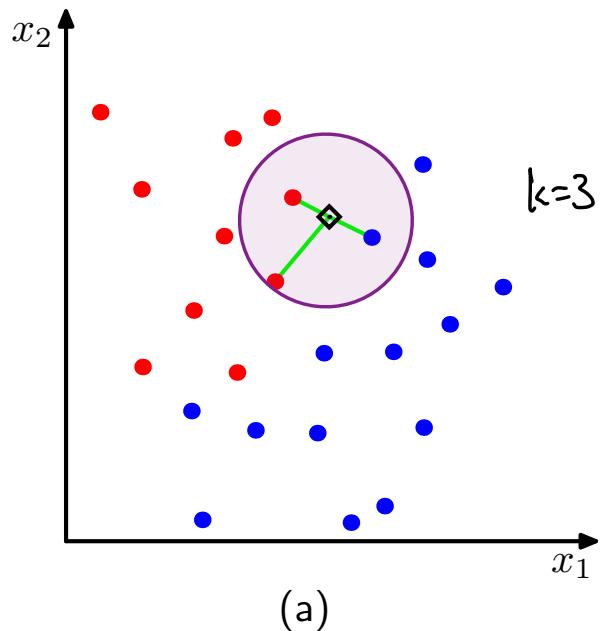
$$\frac{k_n^{(k)}}{k_n} > \frac{k_n^{(j)}}{k_n}$$

$$\therefore \text{if } k_n^{(k)} > k_n^{(j)} \quad \forall j \neq k \Rightarrow \underline{x} \in S_k$$

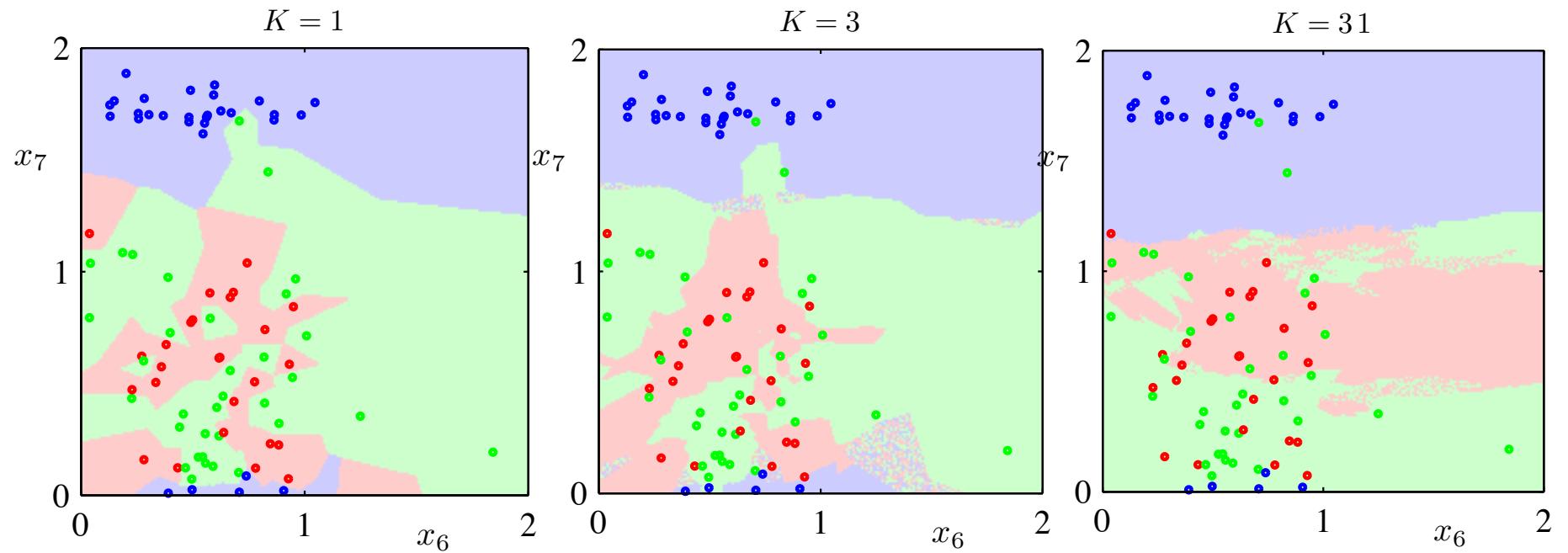
Bayes min.error dec. rule for estimates based on kNN or
L KDE with binary-valued window function.



Ex's - see below

Ex $C = 2$ classes

Bishop Figure 2.27



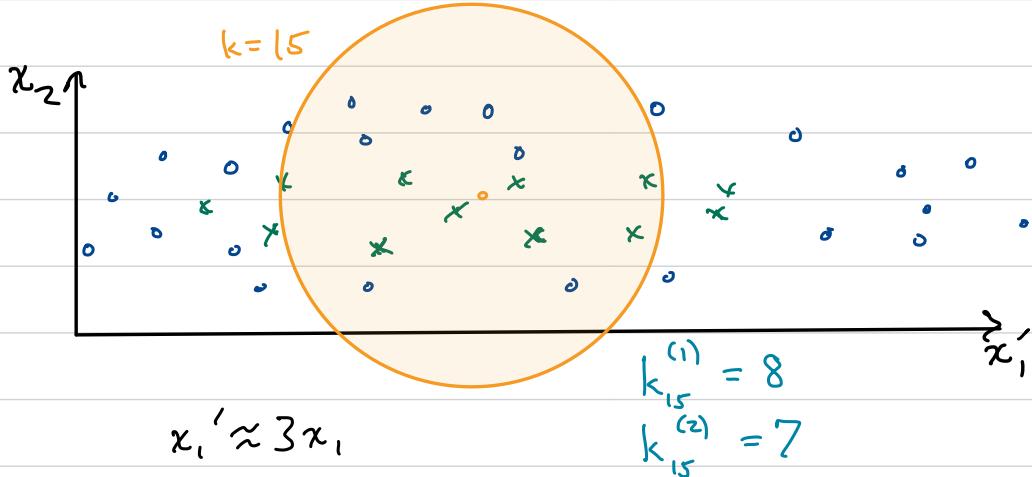
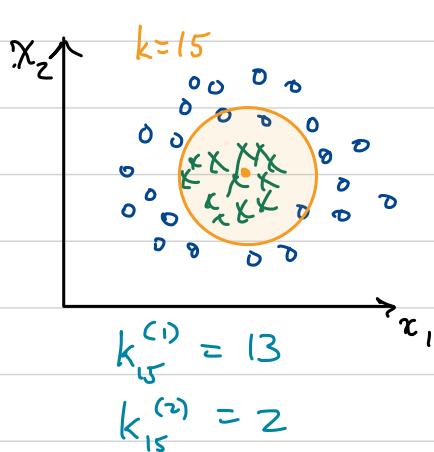
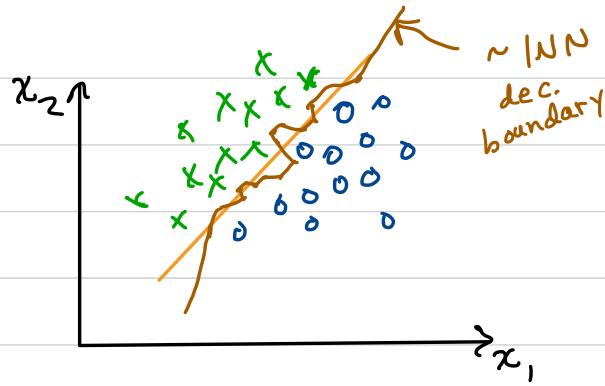
Bishop Figure 2.28

kNN classification on real dataset (oil data (Bishop)), $C=3$.

Comments on kNN and KDE

1. Normalization matters!

- Sensitive to differing scale sizes of different features



2. Choice of k (kNN) or kernel width h (KDE) affects amount of smoothing and resolution (i.e. underfit / overfit).

3. Can be computationally slow for large N or large D .

- There are algorithms to speed up computation (e.g., tree search methods), usually giving an approximate result.