

Statistical Learning I

1. Organization
2. Statistical Learning Intro:
 - a) supervised learning: classification, regression (overfitting)
 - b) simplistic: linear regression for classification (overfitting 2)
 - c) nearest neighbour (overfitting 3)
3. Statistical decision theory:
 - EPE (overfitting 4), Bayes optimal decision
- 5. Bias-variance (overfitting 5)
- 4. Curse of dimensionality ~~not~~ niet aan toegekomen
6. History, applications, beyond classification/regression

1. Organization

Tim van Erven

TA: Kevin Duisters

www.timvanerven.nl/teaching/statlearn2015

(deelnemerslijst uitdelen, boek laten zien)

- 2 large homework assignments > 5.5
- open book exam: Jan 6 > 5.5

HW1: Nov 12 - Dec 1

HW2: Dec 1 - end of Dec

(first two weeks: 3 hours of lectures

then : 2 to 3 hours + work on homework)

Enroll in Blackboard + Usis

②

2a) Supervised Learning

teacher shows us what to do)

Training data $T = \begin{pmatrix} y_1 \\ x_1 \end{pmatrix}, \begin{pmatrix} y_2 \\ x_2 \end{pmatrix}, \dots, \begin{pmatrix} y_N \\ x_N \end{pmatrix}$

x_i : feature vectors

y_i : class / response variable

Goal: learn function $\hat{f}: X \rightarrow Y$
from model \hat{F} .

Evaluate \hat{f} on test data:

new x from same source,

predict corresponding y by $\hat{y} = \hat{f}(x)$.

Assume: $\begin{pmatrix} y_i \\ x_i \end{pmatrix}$ independent from same probability distribution P^* .

Avoid further assumptions on P^* .

(P^* can be very complicated.)

Classification:

$\begin{pmatrix} G_1 \\ x_1 \end{pmatrix}, \dots, \begin{pmatrix} G_N \\ x_N \end{pmatrix}$

G_i is categorical variable (e.g. yes/no or apple/orange/pear)

$G_i \rightarrow Y$

yes $\rightarrow 1$

no $\rightarrow 0$

yes $\rightarrow +1$

no $\rightarrow -1$

apple $\rightarrow 1$

orange $\rightarrow 2$

pear $\rightarrow 3$

apple $\rightarrow \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$

orange $\rightarrow \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$

pear $\rightarrow \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$

"dummy coding"

(see slides)

Linear Regression:

y_i is real-valued (x_i in reals or use dummy coding)
model f can be represented as
linear functions:

$$f_{\beta}(x_i) = \hat{y}_i = \beta_0 + \sum_{j=1}^p x_{i,j} \beta_j$$

$$= x_i^T \beta \quad (x_{i,0} = \text{always } 1)$$

Least Squares: choose $\hat{\beta}$ to minimize

$$RSS(\beta) = \sum_{i=1}^N (y_i - f_{\beta}(x_i))^2 = \sum_{i=1}^N (y_i - x_i^T \beta)^2$$

NB x_i random design, not set by experimenter

(slides)

Important trick: make more features

* Given $x_i \in \mathbb{R}$

* Make $\tilde{x}_{i,0} = 1$

$\tilde{x}_{i,1} = x_i$

$\tilde{x}_{i,2} = x_i^2$

$\tilde{x}_{i,3} = x_i^3$

etc.

Model linear in new features, but not in old features

(polynomial slides)

OVERFITTING is central concern

④ 2b. Linear Regression for Classification

NB. Logistic regression is better!

$G_i \in \{\text{orange, blue}\} \rightarrow y_i \in \{1, 0\}$

$x_i \in \mathbb{R}^2$

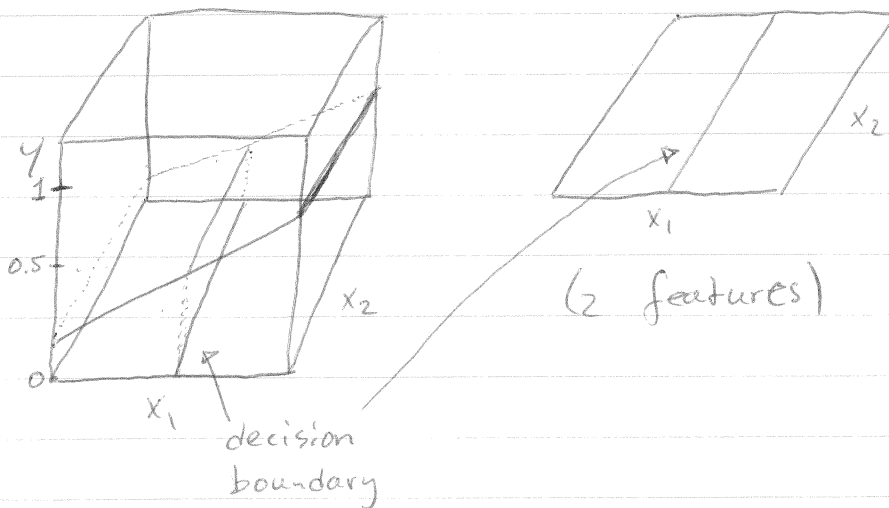
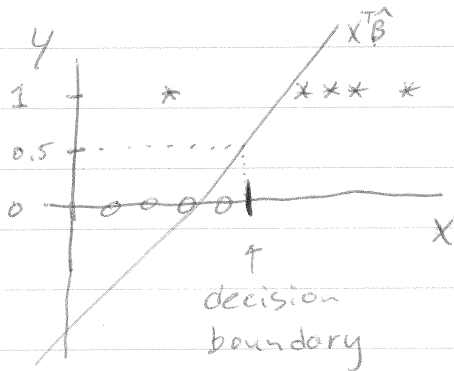
Simple idea:

- use linear regression on this data

- convert back to binary predictions:

$$x^T \hat{\beta} > 0.5 \Rightarrow \hat{f}(x) = 1$$

$$x^T \hat{\beta} \leq 0.5 \Rightarrow \hat{f}(x) = 0$$



(show first figure from book)

- Many errors on training set: avoidable?
- Maybe if we construct more features:

$$x_1, x_2, x_1^2, x_2^2, x_1 \cdot x_2, x_1^2 \cdot x_2, \dots$$

Too few: underfitting
 Too many: overfitting

2C. k-Nearest Neighbour Classification

To classify new x :

- Ask k points in T closest to x
- Choose most common class

For two classes: $y_i \in \{0, 1\}$

$$\hat{y} = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$$

← estimates $Pr(Y=1|X)$ by average in region around x .

$\hat{y} > 0.5 \Rightarrow$ predict 1
 $\hat{y} \leq 0.5 \Rightarrow$ predict 0

k small: very flexible decision boundary, but very unstable estimates based on few neighbours (risks overfitting)

k large: stable estimates, but inflexible decision boundary (risks underfitting)

Model \mathcal{F} is implicit!

If neighbourhoods do not overlap: $\frac{N}{k}$ neighbourhoods

(Show pictures book, what is best k ?)

effective nr. of parameters, because estimate \hat{y} for each neighbourhood.

(Need statistical decision theory!)

(6)

3 Statistical Decision Theory

Learn \hat{f} from training data $(y_1, x_1), \dots, (y_n, x_n) \stackrel{\text{i.i.d.}}{\sim} P^*$
Evaluate on new $(x, y) \sim P^*$

Loss function $L(y, \hat{y})$ measures how bad it is
if truth is y , but we predict $\hat{y} = \hat{f}(x)$
(smaller is better)

Regression: $L(y, \hat{y}) = (\hat{y} - y)^2$ "squared error"

Classification: $L(y, \hat{y}) = \begin{cases} 0 & \text{if } \hat{y} = y \\ 1 & \text{if } \hat{y} \neq y \end{cases}$ "0/1-loss"

Other choices possible?

(depends on what is important in
your application)

$(x, y) \sim P^*$ random, so want small
expected prediction error:

$$\text{EPE}(f) = \mathbb{E}_{x, y \sim P^*} [L(y, f(x))]$$

↑
often called "risk"

Classification:

$$\begin{aligned} \text{EPE}(f) &= \Pr(y \neq f(x)) \cdot 1 + \Pr(y = f(x)) \cdot 0 \\ &= \Pr(y \neq f(x)) \end{aligned}$$

Bayes Optimal Decisions

Bayes optimal f minimizes $\text{EPE}(f)$ over
all possible functions:

$$f_B = \underset{f}{\text{argmin}} \mathbb{E}_{x, y} [L(y, f(x))]$$

$$f_B(x) = \underset{\hat{y}}{\operatorname{argmin}} \mathbb{E}_{P^*(Y|X)} [L(Y, \hat{y})]$$

- unknown, because depends on P^*
- No learning alg can do better (by definition)
- Nothing to do with "Bayesian statistics"

Classification with 0/1-loss:

$$\begin{aligned} f_B(x) &= \underset{\hat{y}}{\operatorname{argmin}} P_{\hat{y}}^*(Y \neq \hat{y} | X) \\ &= \underset{\hat{y}}{\operatorname{argmax}} P_{\hat{y}}^*(Y = \hat{y} | X) \\ &= \begin{cases} 1 & \text{if } P^*(Y=1|X) > 0.5 \\ 0 & \text{if } P^*(Y=1|X) \leq 0.5 \end{cases} \quad (\text{for two classes}) \end{aligned}$$

Regression with Squared Error:

$$\begin{aligned} f_B(x) &= \underset{\hat{y}}{\operatorname{argmin}} \mathbb{E}_{P^*(Y|X)} [(Y - \hat{y})^2] \\ &= \mathbb{E}_{P^*(Y|X)} [Y] = \mathbb{E}[Y|X] \end{aligned}$$

If $Y = f^*(x) + \text{noise}$ ↙ mean 0

$$\begin{aligned} \text{Then } f_B(x) &= \mathbb{E}[Y|X] = f^*(x) + \mathbb{E}[\text{noise}|X] \\ &= f^*(x) \end{aligned}$$

Empirical Risk Minimization

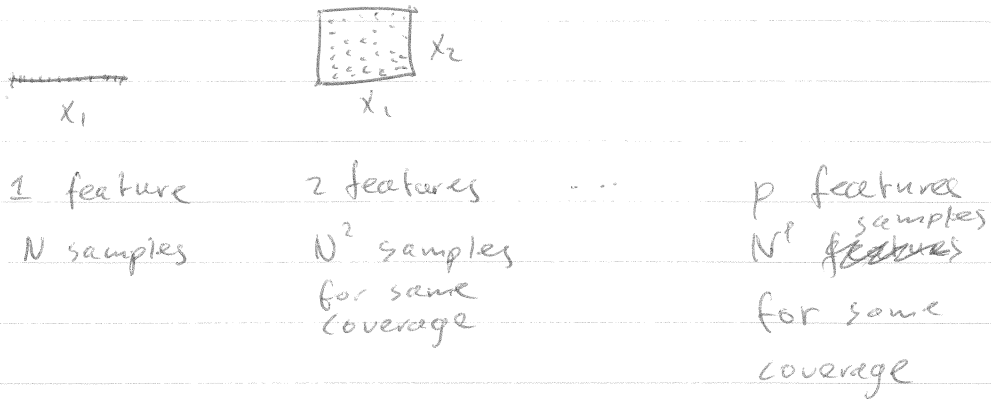
- Approximate P^* by empirical distribution on \mathcal{T} and choose $f \in \mathcal{F}$ that minimizes

$$\frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i))$$

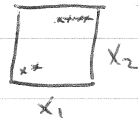
Is this always best? No!

1. Overfitting (see k-NN picture train and test error as function of k)
2. Restricted to model \mathcal{F} .

4. Curse of Dimensionality



Conclusion: In high dimensions you only see samples from very small part of domain.



"A high-dimensional space is a lonely place."
Bernard Schölkopf

5. Bias-variance decomposition for Squared Error

Understand balance between under and overfitting

Thm: Let $\bar{f}(x) = \mathbb{E}_{\mathcal{F}}[\hat{f}(x)]$. Then

$$\begin{aligned} \mathbb{E}_{\mathcal{T}} \text{EPE}(\hat{f}) &= \mathbb{E}_{x,y} \text{Var}(y|x) \quad \leftarrow \text{Bayes error} \\ &+ \mathbb{E}_{x,y} (f_B(x) - \bar{f}(x))^2 \quad \leftarrow \text{bias (typically smaller with larger } \mathcal{F}) \\ &+ \mathbb{E}_{\mathcal{T}} \mathbb{E}_{x,y} (\hat{f}(x) - \bar{f}(x))^2 \quad \leftarrow \text{variance (typically larger with larger } \mathcal{F}) \end{aligned}$$

Proof:

$$E_{T,x,y} [y - \hat{f}(x)]^2 = E_{T,x,y} (y - f_B(x) + f_B(x) - \hat{f}(x))^2$$

(subscripts belangrijk)

$$= E_{x,y} (y - f_B(x))^2$$

$$+ E_{T,x,y} (y - f_B(x))(f_B(x) - \hat{f}(x))$$

$$+ E_{T,x,y} (f_B(x) - \hat{f}(x))^2$$

$$E_{x,y} (y - f_B(x))^2 = E_x (y - E[y|x])^2 = E_x \text{Var}(y|x)$$

$$E_{T,x,y} (y - f_B(x))(f_B(x) - \hat{f}(x))$$

$$= E_{T,x} (E[y|x] - f_B(x))(f_B(x) - \hat{f}(x))$$

$$= 0$$

$$E_{T,x,y} (f_B(x) - \hat{f}(x))^2 = E_{T,x,y} (f_B(x) - \bar{f}(x) + \bar{f}(x) - \hat{f}(x))^2$$

$$= E_x (f_B(x) - \bar{f}(x))^2 \leftarrow \text{bias}$$

$$+ E_{T,x} (f_B(x) - \bar{f}(x))(\bar{f}(x) - \hat{f}(x))$$

$$+ E_{T,x} (\bar{f}(x) - \hat{f}(x))^2 \leftarrow \text{variance}$$

$$E_{T,x} (f_B(x) - \bar{f}(x))(\bar{f}(x) - \hat{f}(x))$$

$$= E_x (f_B(x) - \bar{f}(x))(\bar{f}(x) - E_T[\hat{f}(x)])$$

$$= 0 \quad \square$$