

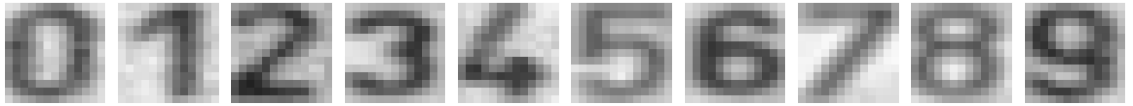
# Classifiers: Naïve Bayes, evaluation

Tomáš Svoboda and Matěj Hoffmann  
thanks to Daniel Novák and Filip Železný, Ondřej Drbohlav

Vision for Robots and Autonomous Systems, Center for Machine Perception  
Department of Cybernetics  
Faculty of Electrical Engineering, Czech Technical University in Prague

May 13, 2020

## Example: Digit recognition/classification



- ▶ **Input:** 8-bit image  $13 \times 13$ , pixel intensities 0 – 255. (0 means black, 255 means white)
- ▶ **Output:** Digit 0 – 9. Decision about the class, classification.
- ▶ **Features:** Pixel intensities ...

---

### Notes

Digit recognition is a very classical example of classification problem. It has been used for decades, and it is used till today, see e.g. [MNIST demo at PyTorch](#)

# Classification as a special case of statistical decision theory

- ▶ Attribute vector  $\vec{x} = (x_1, x_2, \dots)$ : pixels 1, 2, ...
- ▶ **State set  $\mathcal{S}$  = decision set  $\mathcal{D} = \{0, 1, \dots, 9\}$ .**
- ▶ **State = actual class, Decision = recognized class**
- ▶ Loss function:

$$l(s, d) = \begin{cases} 0, & d = s \\ 1, & d \neq s \end{cases}$$

$$\delta^*(\vec{x}) = \arg \min_d \sum_s \underbrace{l(s, d)}_{0 \text{ if } d=s} P(s|\vec{x}) = \arg \min_d \sum_{s \neq d} P(s|\vec{x})$$

Obviously  $\sum_s P(s|\vec{x}) = 1$ , then:

$$P(d|\vec{x}) + \sum_{s \neq d} P(s|\vec{x}) = 1$$

Inserting into above:

$$\delta^*(\vec{x}) = \arg \min_d [1 - P(d|\vec{x})] = \arg \max_d P(d|\vec{x})$$

3 / 24

## Notes

We are using different word – *classification* instead of *decision* but the reasoning and methods can be well applied in both. In classification problem we usually treat all mistakes – wrong classifications – equally painful, contrary to decision problem – remember “What to cook for dinner” problem?

# Bayes classification in practice

- ▶ Usually we are not given  $P(s|\vec{x})$
- ▶ It has to be estimated from already classified examples – training data
- ▶ For discrete  $\vec{x}$ , training examples  $(\vec{x}_1, s_1), (\vec{x}_2, s_2), \dots, (\vec{x}_l, s_l)$ 
  - ▶ so-called i.i.d (independent, identically distributed) multiset
  - ▶ every  $(\vec{x}_i, s_i)$  is drawn independently from  $P(\vec{x}, s)$
- ▶ Without knowing anything about the distribution, a non-parametric estimate:

$$P(s|\vec{x}) = \frac{P(\vec{x}, s)}{P(\vec{x})} \approx \frac{\# \text{ examples where } \vec{x}_i = \vec{x} \text{ and } s_i = s}{\# \text{ examples where } \vec{x}_i = \vec{x}}$$

- ▶ Hard in practice:
  - ▶ To reliably estimate  $P(s|\vec{x})$ , the number of examples grows exponentially with the number of elements of  $\vec{x}$ .
    - ▶ e.g. with the number of pixels in images
    - ▶ curse of dimensionality
    - ▶ denominator often 0

4 / 24

---

## Notes

Why hard? Way too many various  $\vec{x}$ . Think about simple binary  $10 \times 10$  image -  $\vec{x}$  contains 0, 1, position matters.

What is the total number of unique images? Think binary,  $1 \times 8$  binary image?

What is the difference between set and multiset?

Reminder about math notation. In literature, vectors are mostly denoted by bold lower case  $\mathbf{x}$ . In lectures, we use  $\vec{x}$  to match notation used on blackboard. It is difficult to write bold with a chalk.

# Naïve Bayes classification

- ▶ For efficient classification we must thus rely on additional assumptions.
- ▶ In the exceptional case of **statistical independence** between components of  $\vec{x}$  for each class  $s$  it holds

$$P(\vec{x}|s) = P(x[1]|s) \cdot P(x[2]|s) \cdot \dots$$

- ▶ Use simple Bayes law and maximize:

$$P(s|\vec{x}) = \frac{P(\vec{x}|s)P(s)}{P(\vec{x})} = \frac{P(s)}{P(\vec{x})} P(x[1]|s) \cdot P(x[2]|s) \cdot \dots =$$

- ▶ No combinatorial curse in estimating  $P(s)$  and  $P(x[i]|s)$  separately for each  $i$  and  $s$ .
- ▶ No need to estimate  $P(\vec{x})$ . (Why?)
- ▶  $P(s)$  may be provided apriori.
- ▶ **naïve** = when used despite statistical dependence

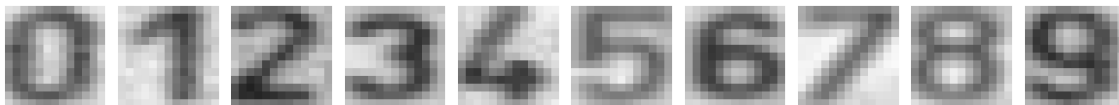
5 / 24

---

## Notes

Why naïve at all? Consider  $N$ - dimensional space,  $8$ - bit values. Instead of problem  $8^N$  we have  $8 \times N$  problem. Think about statistical independence. Example1: person's weight and height. Are they independent? Example2: pixel values in images.

## Example: Digit recognition/classification



- ▶ **Input:** 8-bit image  $13 \times 13$ , pixel intensities 0 – 255. (0 means black, 255 means white)
- ▶ **Output:** Digit 0 – 9. Decision about the class, classification.
- ▶ **Features:** Pixel intensities ...

Collect data , ...

- ▶  $P(\vec{x})$ . What is the dimension of  $\vec{x}$ ? How many possible images?
- ▶ Learn  $P(\vec{x}|s)$  per each class (digit).
- ▶ Classify  $s^* = \operatorname{argmax}_s P(s|\vec{x})$ .

6 / 24

---

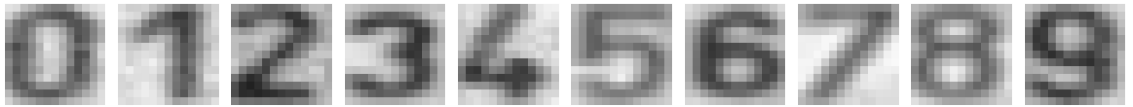
### Notes

We can create many more features than just pixel intensities. But first things first.

We are assuming all errors are equally important - minimizing the number of wrong decisions.

Dimension of  $\vec{x}$  is  $13 \times 13 = 169$ . There are  $255^{169}$  possible images.

## Example: Digit recognition/classification



- ▶ **Input:** 8-bit image  $13 \times 13$ , pixel intensities 0 – 255. (0 means black, 255 means white)
- ▶ **Output:** Digit 0 – 9. Decision about the class, classification.
- ▶ **Features:** Pixel intensities ...

Collect **data** , ...

- ▶  $P(\vec{x})$ . What is the dimension of  $\vec{x}$ ? How many possible images?
- ▶ Learn  $P(\vec{x}|s)$  per each class (digit).
- ▶ Classify  $s^* = \operatorname{argmax}_s P(s|\vec{x})$ .

6 / 24

---

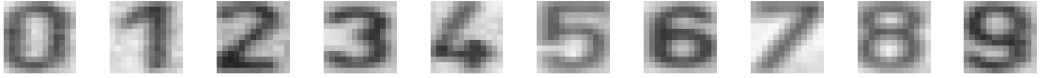
### Notes

We can create many more features than just pixel intensities. But first things first.

We are assuming all errors are equally important - minimizing the number of wrong decisions.

Dimension of  $\vec{x}$  is  $13 \times 13 = 169$ . There are  $255^{169}$  possible images.

From images to  $\vec{x}$

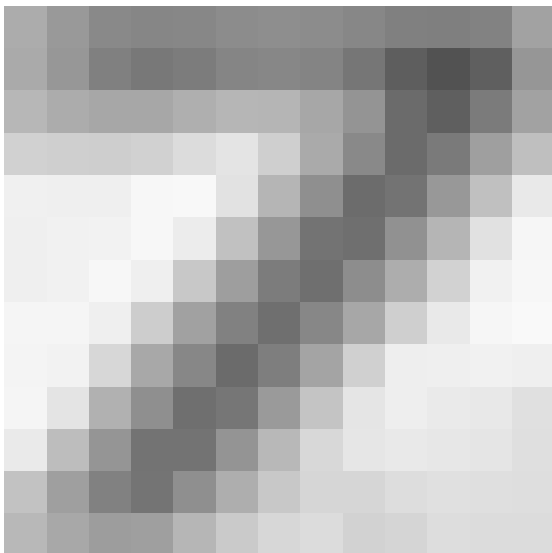


---

Notes



# Conditional probabilities, likelihoods



- ▶ Apriori digit probabilities  $P(s_k)$
- ▶ Likelihoods for pixels.  $P(x_{r,c} = I_i | s_k)$

8 / 24

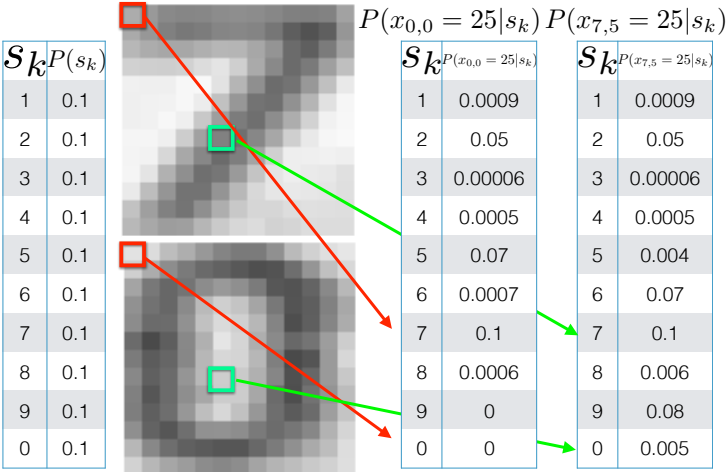
---

## Notes

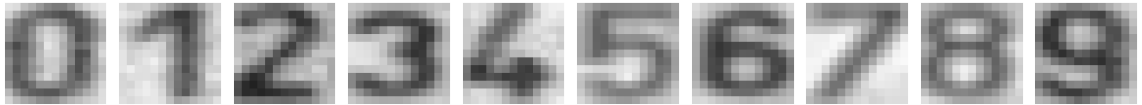
A lexical, especially for Czech speakers. *probability* as well as *likelihood* can be translated as *pravděpodobnost*. I suggest the following mental model than can work for our purposes.

- **Probability** is related to the future events (unknown outcome). E.g. what is the probability of selecting blue box? What is the probability that a random zup number begins with 7?
- **Likelihood** refers to past events (known outcome). In my data, how many images of 7 have dark pixel in top right corner? We can think about relative frequency (relativní četnost).

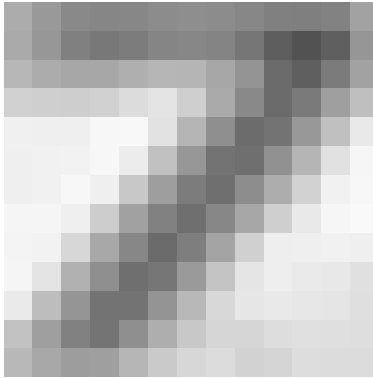
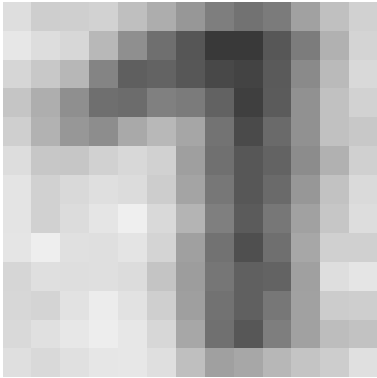
# Conditional likelihoods



## Unseen events



Images  $13 \times 13$ , intensities 0 – 255, 100 exemplars per each class.



10 / 24

---

### Notes

Think about the problem of classifying numerals. Some  $P(x_{r,c} = l \mid s) = 0$ . What about an example:

$$P(x_{0,0} = 100 \mid s = 7) = 0.05$$

$$P(x_{0,0} = 101 \mid s = 7) = 0$$

$$P(x_{0,0} = 102 \mid s = 7) = 0.06$$

A new (not in training) query image with  $x_{0,0} = 101$ . How would you classify?

# Laplace smoothing (“additive smoothing”)

$$P(x) = \frac{\text{count}(x)}{\text{total samples}}$$

Problem:  $\text{count}(x) = 0$

Pretend you see the (any) sample one more time.

$$P_{\text{LAP}}(x) = \frac{c(x) + 1}{\sum_x [c(x) + 1]}$$

$$P_{\text{LAP}}(x) = \frac{c(x) + 1}{N + |X|}$$

where  $N$  is the number of (total) observations;  $|X|$  is the number of possible values  $X$  can take (cardinality).

# Laplace smoothing (“additive smoothing”)

$$P(x) = \frac{\text{count}(x)}{\text{total samples}}$$

Problem:  $\text{count}(x) = 0$

Pretend you see the (any) sample one more time.

$$P_{\text{LAP}}(x) = \frac{c(x) + 1}{\sum_x [c(x) + 1]}$$

$$P_{\text{LAP}}(x) = \frac{c(x) + 1}{N + |X|}$$

where  $N$  is the number of (total) observations;  $|X|$  is the number of possible values  $X$  can take (cardinality).

# Laplace smoothing (“additive smoothing”)

$$P(x) = \frac{\text{count}(x)}{\text{total samples}}$$

Problem:  $\text{count}(x) = 0$

Pretend you see the (any) sample one more time.

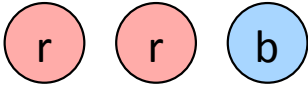
$$P_{\text{LAP}}(x) = \frac{c(x) + 1}{\sum_x [c(x) + 1]}$$

$$P_{\text{LAP}}(x) = \frac{c(x) + 1}{N + |X|}$$

where  $N$  is the number of (total) observations;  $|X|$  is the number of possible values  $X$  can take (cardinality).

$$P_{LAP}(X) = ?$$

Observation:



What is  $P_{LAP}(X = \text{red})$  and  $P_{LAP}(X = \text{blue})$ ?

A:  $P_{LAP}(X = \text{red}) = 7/10$ ,  $P_{LAP}(X = \text{blue}) = 3/10$

B:  $P_{LAP}(X = \text{red}) = 2/3$ ,  $P_{LAP}(X = \text{blue}) = 1/3$

C:  $P_{LAP}(X = \text{red}) = 3/5$ ,  $P_{LAP}(X = \text{blue}) = 2/5$

D: None of the above.

---

Notes



$$P_{ML}(X) =$$

$$P_{LAP}(X) =$$

originally:

- $P(\text{red}) = 2/3$
- $P(\text{blue}) = 1/3$

after Laplace smoothing – adding one red ball and blue ball to the actual observations:

- $P_{LAP}(\text{red}) = (2 + 1)/(2 + 1 + 1 + 1) = 3/5$
- $P_{LAP}(\text{blue}) = (1 + 1)/(2 + 1 + 1 + 1) = 2/5$

this slide: courtesy of P. Abeel, <http://ai.berkeley.edu>. 21st lecture of CS 188.

# Laplace smoothing - as a hyperparameter $k$

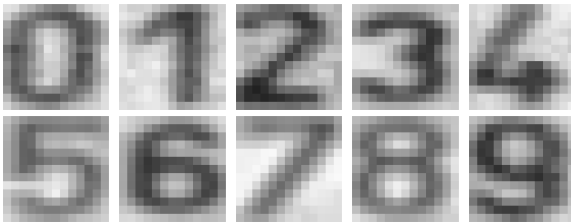
Pretend you see every sample  $k$  extra times:

$$P_{\text{LAP}}(x) = \frac{c(x) + k}{\sum_x [c(x) + k]}$$

$$P_{\text{LAP}}(x) = \frac{c(x) + k}{N + k|X|}$$

For conditional, smooth each condition independently

$$P_{\text{LAP}}(x|s) = \frac{c(x, s) + k}{c(s) + k|X|}$$



What is  $|X|$  equal to?

- A: 10
- B: 2
- C: 256
- D: None of the above

13 / 24

---

## Notes

Hyperparameter would be tuned along with your classifier

For  $k = 100$  and blue and red, you would get:

- $P_{\text{LAP}}(\text{red}) = (2 + 100)/(3 + 100 * 2) = 102/203$
- $P_{\text{LAP}}(\text{blue}) = (1 + 100)/(3 + 100 * 2) = 101/203$

In this case, smoothing ("prior") would dominate over the observations - shifting estimate from empirical to uniform.

In the digit recognition from pixels example: 256 intensity values;  $13 \times 13 = 169$  pixels: Applying Laplace smoothing with  $k = 1$  to  $P(x)$  (prior probability of a particular pixel will take an intensity value  $i$ ):  $P(x_{r,c} = i) = (c(x) + 1)/(N + 256)$

Conditional: relevant for the Naïve Bayes case.



# Laplace smoothing - as a hyperparameter $k$

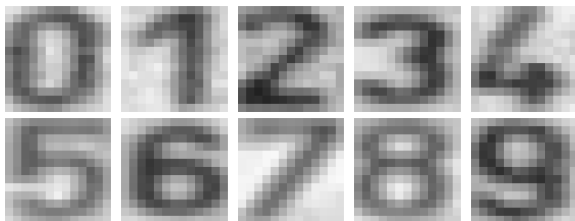
Pretend you see every sample  $k$  extra times:

$$P_{\text{LAP}}(x) = \frac{c(x) + k}{\sum_x [c(x) + k]}$$

$$P_{\text{LAP}}(x) = \frac{c(x) + k}{N + k|X|}$$

For conditional, smooth each condition independently

$$P_{\text{LAP}}(x|s) = \frac{c(x, s) + k}{c(s) + k|X|}$$



What is  $|X|$  equal to?

- A: 10
- B: 2
- C: 256
- D: None of the above

13 / 24

---

## Notes

Hyperparameter would be tuned along with your classifier

For  $k = 100$  and blue and red, you would get:

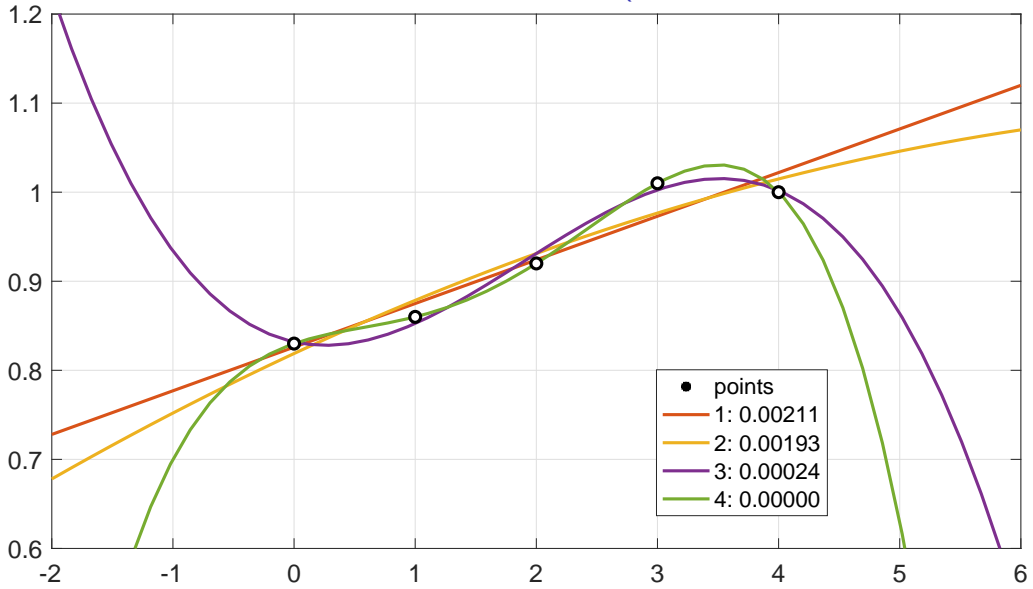
- $P_{\text{LAP}}(\text{red}) = (2 + 100)/(3 + 100 * 2) = 102/203$
- $P_{\text{LAP}}(\text{blue}) = (1 + 100)/(3 + 100 * 2) = 101/203$

In this case, smoothing ("prior") would dominate over the observations - shifting estimate from empirical to uniform.

In the digit recognition from pixels example: 256 intensity values;  $13 \times 13 = 169$  pixels: Applying Laplace smoothing with  $k = 1$  to  $P(x)$  (prior probability of a particular pixel will take an intensity value  $i$ ):  $P(x_{r,c} = i) = (c(x) + 1)/(N + 256)$

Conditional: relevant for the Naïve Bayes case.

# What is the right degree of polynomial (hyperparameter of a regressor)



## Notes

see the `overfit.m` demo

# Generalization and overfitting

- ▶ Data: training, validating, testing . Wanted classifier performs well on what data?
- ▶ Overfitting: too close to training, poor on testing.

# Generalization and overfitting

- ▶ Data: training, validating, testing . Wanted classifier performs well on what data?
- ▶ Overfitting: too close to training, poor on testing.

# Training and testing

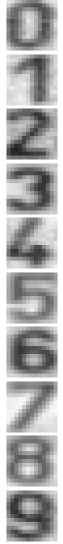
**Data** labeled instances.

- ▶ Training set
- ▶ Held-out (validation) set
- ▶ Testing set.

**Features** : Attribute-value pairs.

**Learning cycle:**

- ▶ **Learn** parameters (e.g. probabilities) on training set.
- ▶ **Tune** hyperparameters on held-out (validation) set.
- ▶ **Evaluate** performance on testing set.



---

## Notes

Training set - biggest part.

# How to evaluate a classifier? Confusion table

Matching table for test set

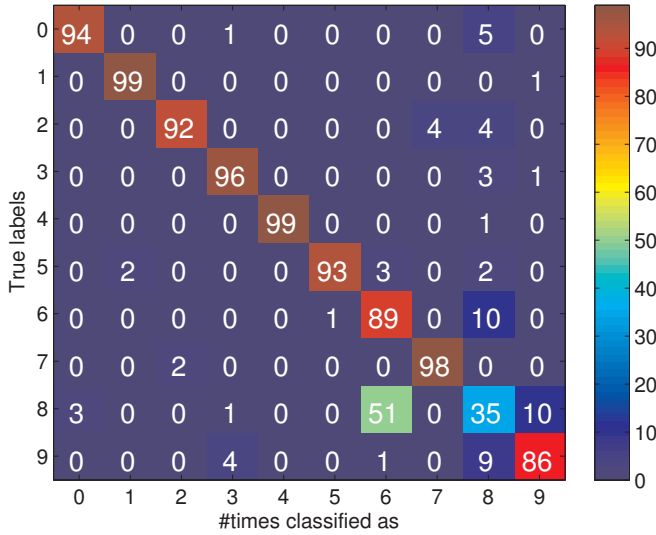


Figure from [5]

## Notes

A result for a one particular classifier and its setting (parameters), one particular testing set

# Precision and Recall, and ...

Consider digit **detection** (is there a digit?) or SPAM/HAM classification.

**Recall** :

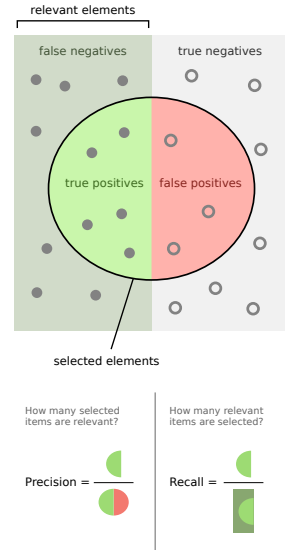
- ▶ How many relevant items are selected?
- ▶ Are we missing some items?
- ▶ Also called: **True positive rate** (TPR), sensitivity, hit rate ...

**Precision**

- ▶ How many selected items are relevant?
- ▶ Also called: Positive predictive value

**False positive rate** (FPR)

- ▶ Probability of false alarm



By Walber - Own work, CC BY-SA 4.0, <https://commons.wikimedia.org/w/index.php?curid=36926283>

18 / 24

## Notes

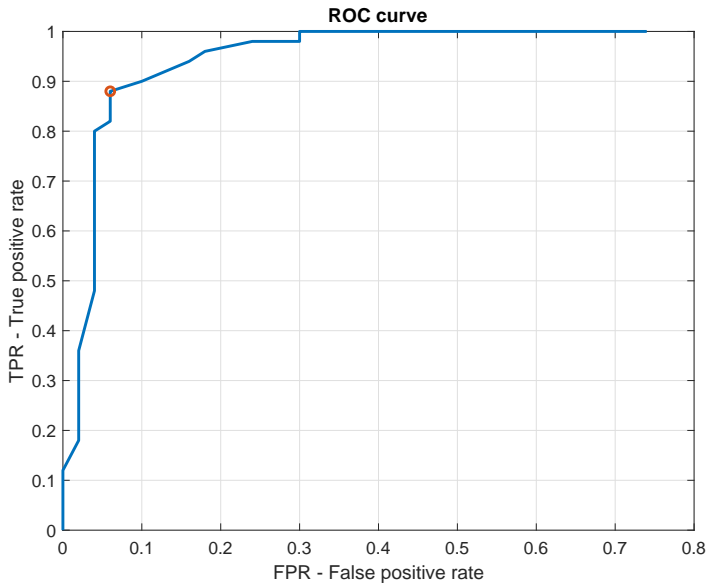
$$\text{TPR} = \frac{\text{TP}}{P} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{FPR} = \frac{\text{FP}}{N} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

Think about TPR vs FPR graph, what is the best classifier?

# ROC – Receiver operating characteristics curve



$$\text{TPR} = \frac{\text{TP}}{P} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$
$$\text{FPR} = \frac{\text{FP}}{N} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

## Notes

- How do you slide along the curve?
- What is the meaning of the diagonal?
- What would be the shape of the curve for the ideal/worst classifier?
- How would you compare various curve and select the best classifier?
- Think/read about other ways to evaluate/visualise classification results.



# Discriminant functions $f(\vec{x}, s)$

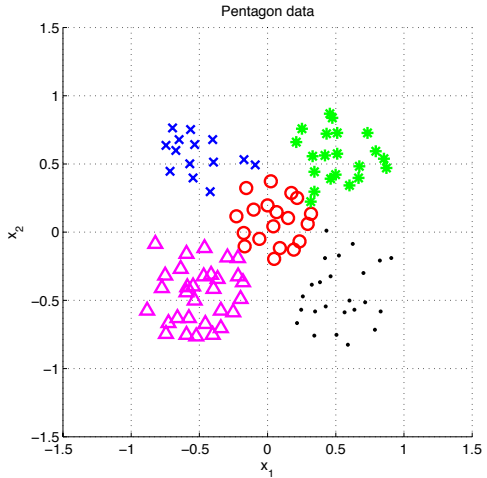
$$s^* = \operatorname{argmax}_{s \in \mathcal{S}} f(\vec{x}, s)$$

Conditional likelihoods:  $\mathcal{N}(\vec{x} | \vec{\mu}_s, \Sigma_s)$

$$\frac{1}{2\pi |\Sigma_s|^{1/2}} \exp\left\{-\frac{1}{2}(\vec{x} - \vec{\mu}_s)^\top \Sigma_s^{-1}(\vec{x} - \vec{\mu}_s)\right\}$$

Bayes:

$$s^* = \operatorname{argmax}_{s \in \mathcal{S}} P(s | \vec{x}) = \frac{P(\vec{x} | s)P(s)}{P(\vec{x})}$$



$$f(\vec{x}, s) = P(s) \frac{1}{2\pi |\Sigma_s|^{1/2}} \exp\left\{-\frac{1}{2}(\vec{x} - \vec{\mu}_s)^\top \Sigma_s^{-1}(\vec{x} - \vec{\mu}_s)\right\}$$

20 / 24

## Notes

Normal distribution for general dimensionality D:

$$\mathcal{N}(\vec{x} | \vec{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\vec{x} - \vec{\mu})^\top \Sigma^{-1}(\vec{x} - \vec{\mu})\right\}$$

How about learning  $f(\vec{x}, s)$  directly without explicit modeling of underlying probabilities?

What about  $f(\vec{x}, s) = \vec{w}_s^\top \vec{x} + w_{s0}$

## Towards linear classifier, geometrical thoughts ...

$$f(\vec{x}, s) = P(s) \frac{1}{2\pi |\Sigma_s|^{1/2}} \exp\left\{-\frac{1}{2}(\vec{x} - \vec{\mu}_s)^\top \Sigma_s^{-1}(\vec{x} - \vec{\mu}_s)\right\}$$

# Product of many small numbers ...

$$P(s|\vec{x}) = \frac{P(\vec{x}|s)P(s)}{P(\vec{x})} = \frac{P(s)}{P(\vec{x})} P(x[1]|s) \cdot P(x[2]|s) \cdot \dots$$

$P(\vec{x})$  not needed, .....

$$\log(P(x[1]|s)P(x[2]|s) \dots) = \log(P(x[1]|s)) + \log(P(x[2]|s)) + \dots$$

---

## Notes

just try

- `prod(rand(1,100))` and `prod(rand(1,10000))` in Matlab.
- `prod(rand(1,100)) == 0` and `prod(rand(1,10000)) == 0` in Matlab.

or in python console:

- `>>> import numpy as np`
- `>>> np.prod(np.random.rand(100))==0`
- `>>> np.prod(np.random.rand(1000))==0`
- `>>> a = np.random.rand(1000)`  
`>>> b = np.random.rand(1000)`  
`>>> np.prod(a)>np.prod(b)`  
False  
`>>> np.prod(a)<np.prod(b)`  
False  
`>>> np.sum(np.log(a))>np.sum(np.log(b))`  
True

Hitting the limit of number representation.

What is the way out?

$P(\vec{x})$  not needed – does not depend on the class.

Laws of logarithms...

# Product of many small numbers ...

$$P(s|\vec{x}) = \frac{P(\vec{x}|s)P(s)}{P(\vec{x})} = \frac{P(s)}{P(\vec{x})} P(x[1]|s) \cdot P(x[2]|s) \cdot \dots$$

$P(\vec{x})$  not needed, .....

$$\log(P(x[1]|s)P(x[2]|s) \dots) = \log(P(x[1]|s)) + \log(P(x[2]|s)) + \dots$$

---

## Notes

just try

- `prod(rand(1,100))` and `prod(rand(1,10000))` in Matlab.
- `prod(rand(1,100)) == 0` and `prod(rand(1,10000)) == 0` in Matlab.

or in python console:

- `>>> import numpy as np`
- `>>> np.prod(np.random.rand(100))==0`
- `>>> np.prod(np.random.rand(1000))==0`
- `>>> a = np.random.rand(1000)`  
`>>> b = np.random.rand(1000)`  
`>>> np.prod(a)>np.prod(b)`  
False  
`>>> np.prod(a)<np.prod(b)`  
False  
`>>> np.sum(np.log(a))>np.sum(np.log(b))`  
True

Hitting the limit of number representation.

What is the way out?

$P(\vec{x})$  not needed – does not depend on the class.

Laws of logarithms...

# References I

Further reading: Chapter 13 and 14 of [4]. Books [1] and [2] are classical textbooks in the field of pattern recognition and machine learning. This lecture has been also inspired by the 21st lecture of CS 188 at <http://ai.berkeley.edu> (e.g., Laplace smoothing). Many Matlab figures created with the help of [3].

[1] Christopher M. Bishop.

*Pattern Recognition and Machine Learning.*

Springer Science+Bussiness Media, New York, NY, 2006.

PDF freely downloadable.

[2] Richard O. Duda, Peter E. Hart, and David G. Stork.

*Pattern Classification.*

John Wiley & Sons, 2nd edition, 2001.

## References II

- [3] Votjěch Franc and Václav Hlaváč.  
Statistical pattern recognition toolbox.  
<http://cmp.felk.cvut.cz/cmp/software/stprtool/index.html>.
  
- [4] Stuart Russell and Peter Norvig.  
*Artificial Intelligence: A Modern Approach*.  
Prentice Hall, 3rd edition, 2010.  
<http://aima.cs.berkeley.edu/>.
  
- [5] Tomáš Svoboda, Jan Kybic, and Hlaváč Václav.  
*Image Processing, Analysis and Machine Vision — A MATLAB Companion*.  
Thomson, Toronto, Canada, 1<sup>st</sup> edition, September 2007.  
<http://visionbook.felk.cvut.cz/>.