

Introduction to Machine Learning

NPFL 054

<http://ufal.mff.cuni.cz/course/npfl054>

Barbora Hladká
hladka@ufal.mff.cuni.cz

Martin Holub
holub@ufal.mff.cuni.cz

Charles University,
Faculty of Mathematics and Physics,
Institute of Formal and Applied Linguistics

Outline

- **Logistic regression**
- **Evaluation of binary classifiers**

Binary classification

Decision boundary

A task of binary classification: $Y = \{0, 1\}$

Decision boundary takes a form of function f and partitions a feature space into two sets, one for each class.

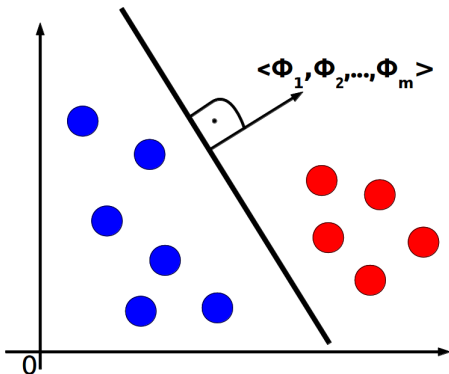
Binary classification

Hyperplane

Hyperplane is a linear decision boundary of the form

$$\Theta^T \mathbf{x} = 0$$

where direction of $\langle \theta_1, \theta_2, \dots, \theta_m \rangle$ is perpendicular to the hyperplane and θ_0 determines position of the hyperplane with respect to the origin



Binary classification

Hyperplane

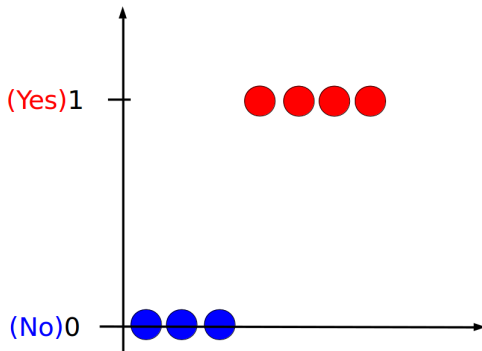
- point if $m = 1$, line if $m = 2$, plane if $m = 3$, ...
- we can use hyperplane for classification so that

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \theta_0 + \theta_1 x_1 + \dots + \theta_m x_m \geq 0 \\ 0 & \text{if } \theta_0 + \theta_1 x_1 + \dots + \theta_m x_m < 0 \end{cases}$$

- **linear classifiers** classify examples using hyperplanes

Binary classification

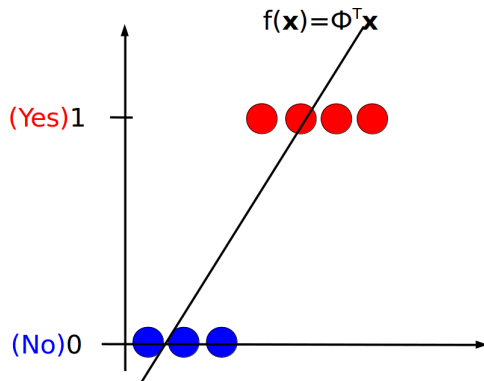
Can we use linear regression?



Binary classification

Can we use linear regression?

Fit the data with a linear function f

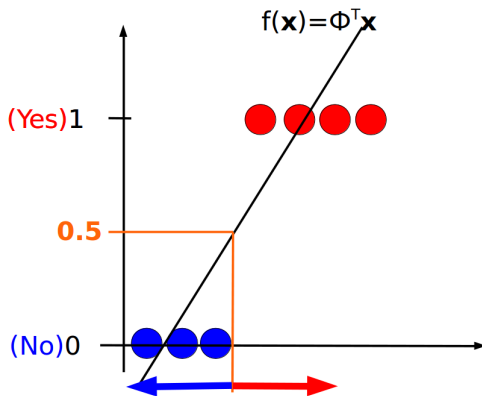


Binary classification

Can we use linear regression?

Classify

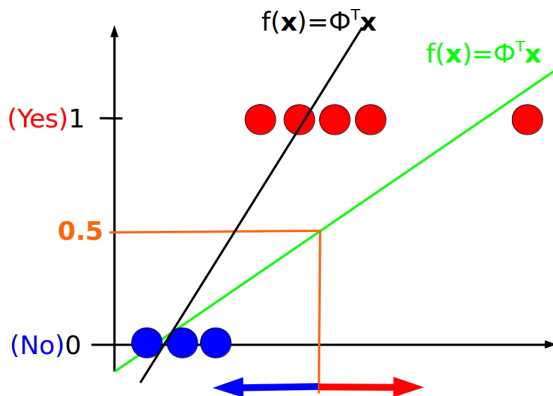
- if $f(\mathbf{x}) \geq 0.5$, predict **1**
- if $f(\mathbf{x}) < 0.5$, predict **0**



Binary classification

Can we use linear regression?

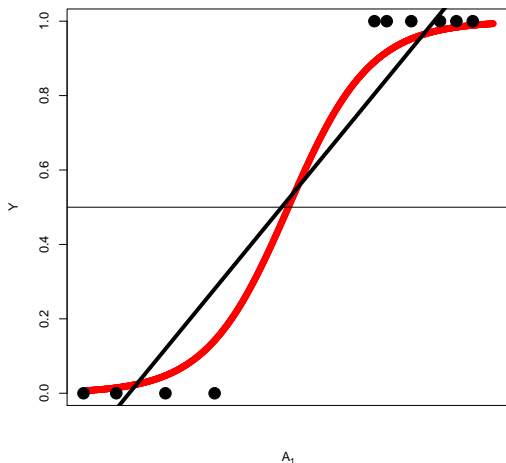
Add one more training instance



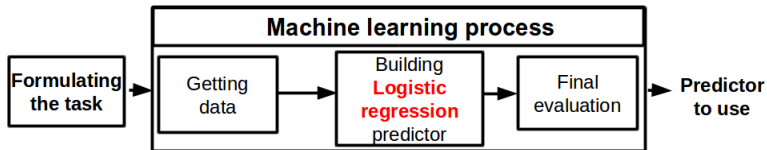
Binary classification

Can we use linear regression?

We are heading for the logistic regression algorithm.



Logistic regression

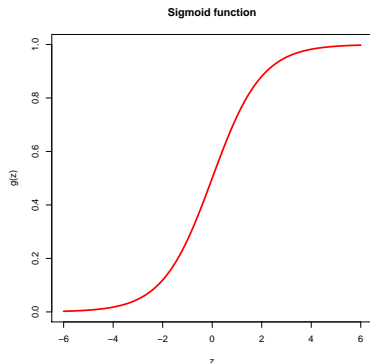


Logistic regression

Logistic regression is a classification algorithm.

Its target hypothesis f for a binary classification has a form of **sigmoid function**

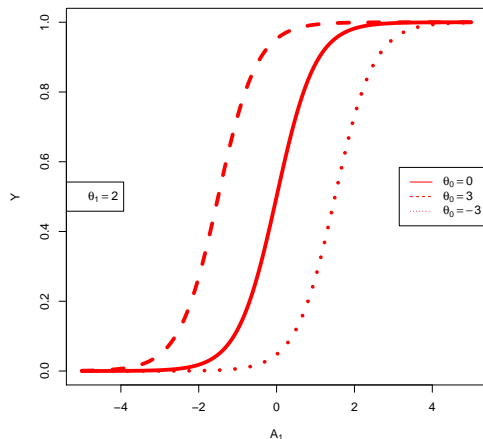
$$f(\mathbf{x}; \Theta) = \frac{1}{1 + e^{-\Theta^T \mathbf{x}}} = \frac{e^{\Theta^T \mathbf{x}}}{1 + e^{\Theta^T \mathbf{x}}}$$



- $g(z) = \frac{1}{1+e^{-z}}$
- $\lim_{z \rightarrow +\infty} g(z) = 1$
- $\lim_{z \rightarrow -\infty} g(z) = 0$

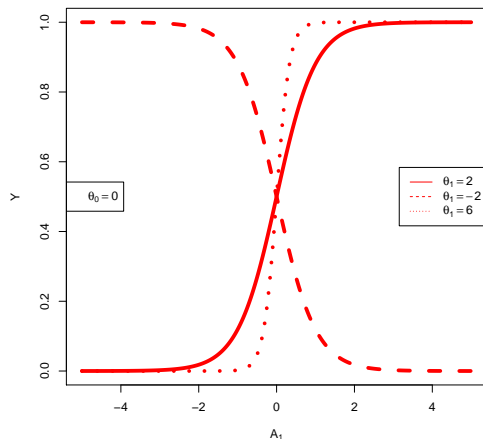
Logistic regression

$$f(\mathbf{x}; \Theta) = \frac{1}{1 + e^{-\theta_0 - \theta_1 x_1}}$$



Logistic regression

$$f(\mathbf{x}; \Theta) = \frac{1}{1 + e^{-\theta_0 - \theta_1 x_1}}$$

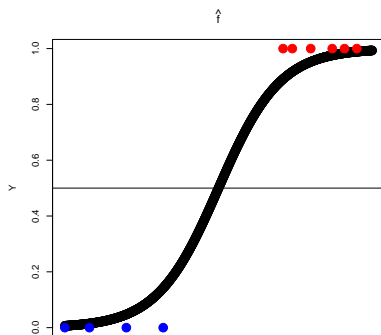


Logistic regression

Classification rule

Predict a target value using $\hat{f}(\mathbf{x}; \hat{\Theta})$ so that

- if $\hat{f}(\mathbf{x}; \hat{\Theta}) \geq 0.5$, i.e. $\hat{\Theta}^\top \mathbf{x} \geq 0$, predict **1**
- if $\hat{f}(\mathbf{x}; \hat{\Theta}) < 0.5$, i.e. $\hat{\Theta}^\top \mathbf{x} < 0$, predict **0**



Logistic regression

Derivation

Interpretation of $f(\mathbf{x}; \Theta)$: it models the conditional probability $\Pr(y = 1|\mathbf{x}; \Theta)$

$$f(\mathbf{x}; \Theta) = \Pr(y = 1|\mathbf{x}; \Theta)$$

1. categorical attribute $Y = \{0, 1\}$
2. ~~$y = \theta_0 + \theta_1 x_1 \cdots + \theta_m x_m$, see above \rightarrow model $\Pr(Y = y|\mathbf{x})$, e.g. $\Pr(Y = 1|\mathbf{x})$~~
3. ~~$\Pr(Y = 1|\mathbf{x}) = \theta_0 + \theta_1 x_1 \cdots + \theta_m x_m$, see above~~
4. Model odds($\Pr(Y = 1|\mathbf{x})$) = $\frac{\Pr(Y=1|\mathbf{x})}{\Pr(Y=0|\mathbf{x})} = \frac{\Pr(Y=1|\mathbf{x})}{1-\Pr(Y=1|\mathbf{x})} \in \langle 0, +\infty \rangle$

Odds, odds ratio

odds = $\Pr(\text{success}) / \Pr(\text{failure})$

Example: Titanic data set

```
> d <- read.csv("train.csv")
> attach(d)
> table(Sex, Survived)
      Survived
Sex      0    1
female  81 233
male   468 109
> detach()
```

- the odds of surviving for male:
 $\Pr(\text{Survived} = 1 | \text{Sex} = \text{male}) / \Pr(\text{Survived} = 0 | \text{Sex} = \text{male}) = \frac{109}{468} = 0.23$
- the odds of surviving for female:
 $\Pr(\text{Survived} = 1 | \text{Sex} = \text{female}) / \Pr(\text{Survived} = 0 | \text{Sex} = \text{female}) = \frac{233}{81} = 2.88$
- the ratio of the odds for female to the odds for male $2.88 / 0.23 = 12.52$

5. Transform $(0, +\infty)$ to $(-\infty, +\infty)$: model

$$\text{logit}(\Pr(Y = 1|\mathbf{x})) = \ln(\text{odds}(\Pr(Y = 1|\mathbf{x}))) = \ln\left(\frac{\Pr(Y = 1|\mathbf{x})}{1 - \Pr(Y = 1|\mathbf{x})}\right)$$

6. Use linear regression

$$\ln\left(\frac{\Pr(Y = 1|\mathbf{x})}{1 - \Pr(Y = 1|\mathbf{x})}\right) = \theta_0 + \theta_1 x_1 + \dots + \theta_m x_m$$

i.e.,

$$\Pr(Y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\theta_0 - \theta_1 x_1 - \dots - \theta_m x_m}}$$

$$f(\mathbf{x}_i; \Theta) = \Pr(Y_i = 1|\mathbf{x}_i; \Theta) = \frac{1}{1 + e^{-\Theta^\top \mathbf{x}_i}}$$

Parameter interpretation

Binary features

- Use female = {1, 0} instead of Sex = {female, male}
- in Linear regression $y = \theta_0 + \theta_1 * \text{female}$
 - θ_0 is the average y for male
 - $\theta_0 + \theta_1$ is the average y for female
 - θ_1 is the average difference in y between female and male
- in Logistic regression $p = \Pr(\text{Survive} = 1 | \mathbf{x}, \Theta)$, $\ln \frac{p}{1-p} = \theta_0 + \theta_1 * \text{female}$
 - **If** female == 0
 - $p = p_1 \rightarrow \ln\left(\frac{p_1}{1-p_1}\right) = \theta_0 \rightarrow \frac{p_1}{1-p_1} = e^{\theta_0}$
 - the intercept θ_0 is the log odds for men
 - **If** female == 1
 - $p = p_2 \rightarrow \frac{p_2}{1-p_2} = e^{\theta_0 + \theta_1}$
 - odds ratio = $\frac{p_2}{1-p_2} / \frac{p_1}{1-p_1} = e^{\theta_1}$
 - the parameter θ_1 is the log odds ratio between female and male

Parameter interpretation

Numerical features

- θ_i gives an average change in $\text{logit}(f(\mathbf{x}))$ with one-unit change in A_i holding all other features fixed

- **Loss function**

$$L(\Theta) = - \sum_{i=1}^n y_i \log P(y_i|\mathbf{x}_i; \Theta) + (1 - y_i) \log(1 - P(y_i|\mathbf{x}_i; \Theta))$$

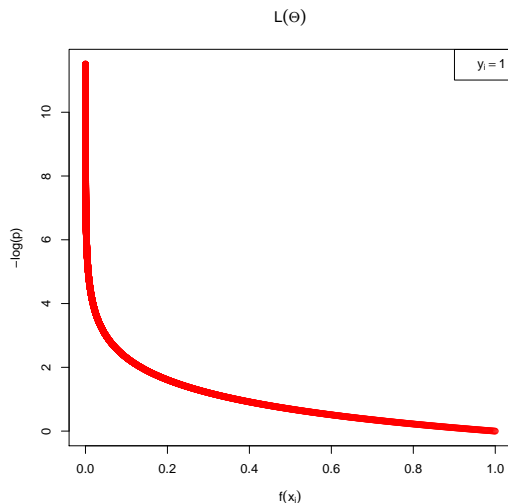
See Maximum Likelihood Principle for derivation of this loss function.

- **Optimization problem**

$$\Theta^* = \operatorname{argmin}_{\Theta} L(\Theta)$$

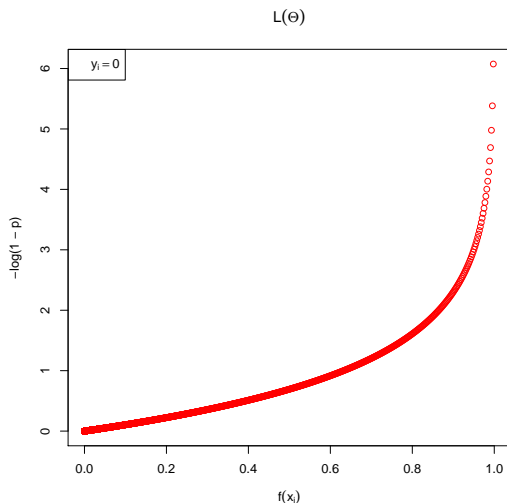
Parameter estimates

$$L(\Theta) = - \sum_{i=1}^n y_i \log P(y_i | \mathbf{x}_i; \Theta) + (1 - y_i) \log(1 - P(y_i | \mathbf{x}_i; \Theta))$$



Parameter estimates

$$L(\Theta) = - \sum_{i=1}^n y_i \log P(y_i | \mathbf{x}_i; \Theta) + (1 - y_i) \log(1 - P(y_i | \mathbf{x}_i; \Theta))$$



Parameter estimates

Gradient Descent Algorithm

repeat until convergence {

$$\Theta^{K+1} := \Theta^K - \alpha \nabla f(\Theta^K)$$

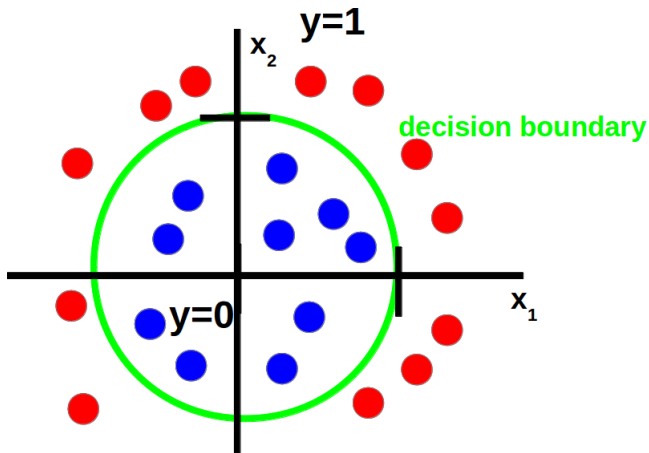
}

– α is a positive step-size hyperparameter

I.e. simultaneously update θ_j , $j = 1, \dots, m$

$$\theta_j^{K+1} := \theta_j^K - \alpha \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i; \Theta^K) - y_i) x_{ij}$$

Non-linear decision boundary



Non-linear decision boundary

- For hyperplane: $f(\mathbf{x}) = g(\theta_0 + \theta_1 x_1 + \dots + \theta_m x_m)$
- Let $f(\mathbf{x}) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$ (a higher degree polynomial)
- Assume $\theta_0 = -1, \theta_1 = 0, \theta_2 = 0, \theta_3 = 1, \theta_4 = 1$
- Predict $y = 1$ if $-1 + x_1^2 + x_2^2 \geq 0$, i.e. $x_1^2 + x_2^2 \geq 1$

Logistic regression

Summary

Classification of \mathbf{x} by \hat{f}^*

- 1 Project \mathbf{x} onto $\hat{\Theta}^*$ to convert it into a real number z in the range $(-\infty, +\infty)$
 - i.e. $z = \hat{\Theta}^{*\top} \mathbf{x}$
- 2 Map z to the range $\langle 0, 1 \rangle$ using the sigmoid function $g(z) = 1/(1 + e^{-z})$
- 3 Classify \mathbf{x} using a classification rule

Multi-class classification

$$|Y| = N, N \geq 3$$

- **One-to-all**

- train N predictors f_k for the pair k -th class and $\{1, \dots, N\} \setminus \{k\}$ classes
- classify \mathbf{x} into the class $k^* = \operatorname{argmax}_k f_k(\mathbf{x})$

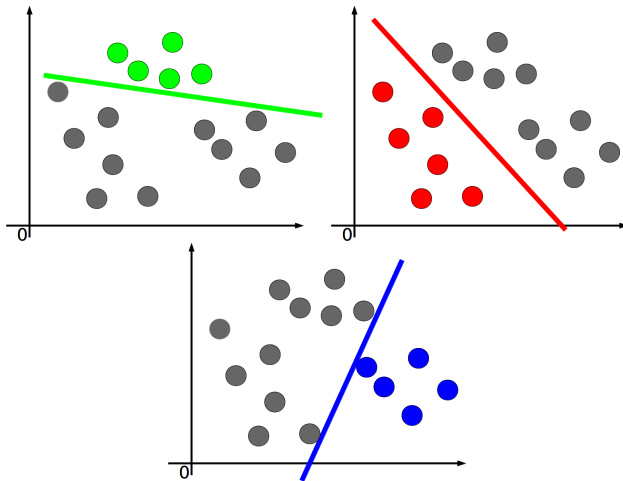
- **One-to-one**

- train $\binom{N}{2}$ classifiers f_i
- classify \mathbf{x} into the class $k^* = \max_{k=1, \dots, N} \sum_{i=1}^{\binom{N}{2}} \delta(f_i(\mathbf{x}) = k)$

Logistic regression

Multi-class classification

One-to-all



Evaluation of binary classifiers

Confusion matrix

Confusion matrix is a square matrix indexed by all possible target class values.

Task: Assign the correct sense of the word *line* in a sentence.

```
** Comparing the predicted values with the true senses **
```

	Prediction					
Truth	cord	division	formation	phone	product	text
cord	268	3	10	7	9	6
division	3	280	1	2	5	3
formation	13	3	225	4	19	4
phone	25	5	2	293	12	10
product	51	10	39	32	1442	72
text	12	1	7	4	28	262

Correctly predicted examples are displayed on the diagonal.

Evaluation of binary classifiers

Confusion matrix

In binary classification tasks examples are sometimes regarded as divided into two disjoint subsets:

- **positive examples** – “to be retrieved” (ones)
- **negative examples** – “not to be retrieved” (zeros)

```
### Example confusion matrix for binary classification
> table(test.true, test.pred)
      prediction
      0      1
true  0 580  69
      1  37 144
>
```

Evaluation of binary classifiers

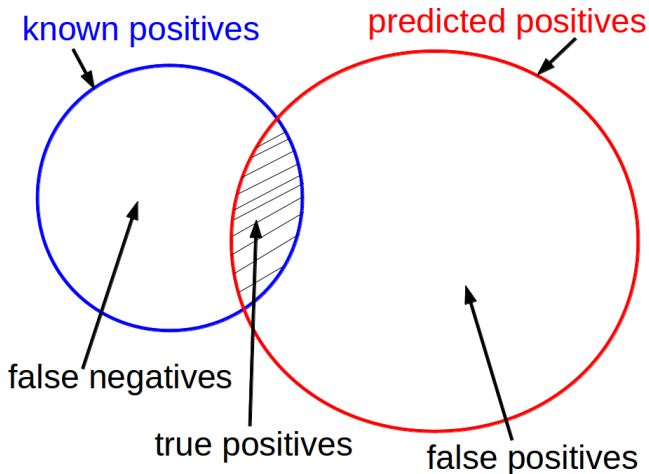
Confusion matrix

		Predicted class	
		Positive	Negative
True class	Positive	True Positive (TP)	False Negative (FN)
	Negative	False Positive (FP)	True Negative (TN)

Explanation

- **'Trues'** are examples correctly classified
- **'Falses'** are examples incorrectly classified
- **'Positives'** were predicted as positives (correctly or incorrectly)
- **'Negatives'** were predicted as negatives (correctly or incorrectly)

Proportion of correctly predicted test examples



Basic performance measures

Measure	Formula
Precision	$TP/(TP+FP)$
Recall/Sensitivity	$TP/(TP+FN)$
Specificity	$TN/(TN+FP)$
Accuracy	$(TP+TN)/(TP+FP+TN+FN)$

Very often you need to **combine both good precision and good recall**. Then you usually use **balanced F-score**, so called **F-measure**

$$F = 2 \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

Summary of Examination Requirements

- Decision boundary, classification rule
- Logistic regression, sigmoid function, probabilistic formulation
- Parameter interpretation
- Multi-class classification
- Confusion matrix for binary classification
- Basic performance measures