

LECTURE 2

Supervised text classification

SUPERVISED TEXT CLASSIFICATION

- Objective
 - Given a set of annotated texts, *i.e.* examples
 - $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\} \quad \forall (x_i, y_i)_{1 \leq i \leq n} : x_i \in \mathcal{X}, y_i \in \mathcal{Y}$
 - Find a function (the decision rule) that predicts the class of a document
 - $f : \mathcal{X} \rightarrow \mathcal{Y}$
- Lecture plan
 - Two parametric models
 - Multinomial naïve Bayse classifier
 - Logistic regression

MULTINOMIAL NAÏVE BAYES CLASSIFIER

MODEL

- Data: a set of annotated documents

- $D = \{(x_1, y_1), \dots, (x_n, y_n)\} \quad \forall (x_i, y_i)_{1 \leq i \leq n} : x_i \in \mathbb{N}^d, y_i \in \{1, \dots, q\}$

- Joint distribution

- $P(X = x, Y = y) = P(Y = y)P(X = x | Y = y)$

- Naïve independence assumption

- $P(X = x | Y = y) = \prod_{j=1}^d P(X_j = x_j | Y = y)$

- $P(X = x, Y = y) = P(Y = y) \prod_{j=1}^d P(X_j = x_j | Y = y)$

MODEL

► Decision rule

$$\text{► } \hat{y} = \operatorname{argmax}_{y \in \{1, 2, \dots, q\}} \left(P(Y = y | X = z) \right)$$

► By Bayes rule we have

$$\text{► } \hat{y} = \operatorname{argmax}_{y \in \{1, 2, \dots, q\}} \left(\frac{P(X = x | Y = y)P(Y = y)}{P(X = x)} \right)$$

$$\text{► } \hat{y} = \operatorname{argmax}_{y \in \{1, 2, \dots, q\}} \left(P(X = x | Y = y)P(Y = y) \right)$$

► Introducing parameters we get

$$\text{► } \hat{y} = \operatorname{argmax}_{y \in \{1, 2, \dots, q\}} \left(g(y) \prod_{j=1}^d g_j(1 | y)^{z_j} \right)$$

PARAMETER ESTIMATION

➤ Maximum likelihood estimation (MLE)

$$\text{➤ } L(D; g, g_j) = \prod_{i=1}^n P(X = x_i, Y = y_i)$$

$$\begin{aligned} \text{➤ } \ell(D; g, g_j) &= \log \left(\prod_{i=1}^n P(X = x_i, Y = y_i) \right) \\ &= \log \left(\prod_{i=1}^n \left(g(y_i) \prod_{j=1}^d g_j(x_{ij} | y_i) \right) \right) \\ &= \sum_{i=1}^n \log \left(g(y_i) \prod_{j=1}^d g_j(x_{ij} | y_i) \right) \\ &= \sum_{i=1}^n \left[\log(g(y_i)) + \sum_{j=1}^d \log g_j(x_{ij} | y_i) \right] \end{aligned}$$

PARAMETER ESTIMATION

- Maximum likelihood estimation (MLE)

- $g(y)$ is the relative frequency of documents of class y in D

- $g(y) = \frac{\sum_{i=1}^n \delta_{y_i=y}}{n}$

- $g_j(1 | y)$ is the relative frequency of word j in documents belonging to class y

- $g_j(1 | y) = \frac{\sum_{i=1}^n \delta_{y_i=y} x_{ij}}{\sum_{i=1}^n \sum_{k=1}^d \delta_{y_i=y} x_{ik}}$

LAPLACE SMOOTHING

- MLE leads to zero estimates for some words, which is problematic given the decision rule

- $\hat{y} = \operatorname{argmax}_{y \in \{1, 2, \dots, q\}} \left(g(y) \prod_{j=1}^d g_j(1 | y)^{z_j} \right)$

- A usual workaround is applying Laplace smoothing
 - Simply increment the frequency of each word in each doc

- $$g_j(1 | y) = \frac{\sum_{i=1}^n (\delta_{y_i=y} x_{ij} + 1)}{\sum_{i=1}^n \sum_{k=1}^d (\delta_{y_i=y} x_{ik} + 1)}$$

LOGISTIC REGRESSION

MODEL

- Data: a set of annotated documents

- $D = \{(x_1, y_1), \dots, (x_n, y_n)\} \quad \forall (x_i, y_i)_{1 \leq i \leq n} : x_i \in \mathbb{R}^d, y_i \in \{0, 1\}$

- Conditional distribution

- $P(Y = 1 | X = x) = \frac{1}{1 + e^{-f(x)}}$

- We recognize the standard sigmoid function

- $S : \mathbb{R} \rightarrow [0; 1] \quad \alpha \mapsto \frac{1}{1 + e^{-\alpha}}$

- Applied to a linear function of x , parameterized by w

- $f : \mathbb{R}^d \rightarrow \mathbb{R}$

- $x \mapsto w_0 + w_1x_1 + w_2x_2 + \dots + w_dx_d$

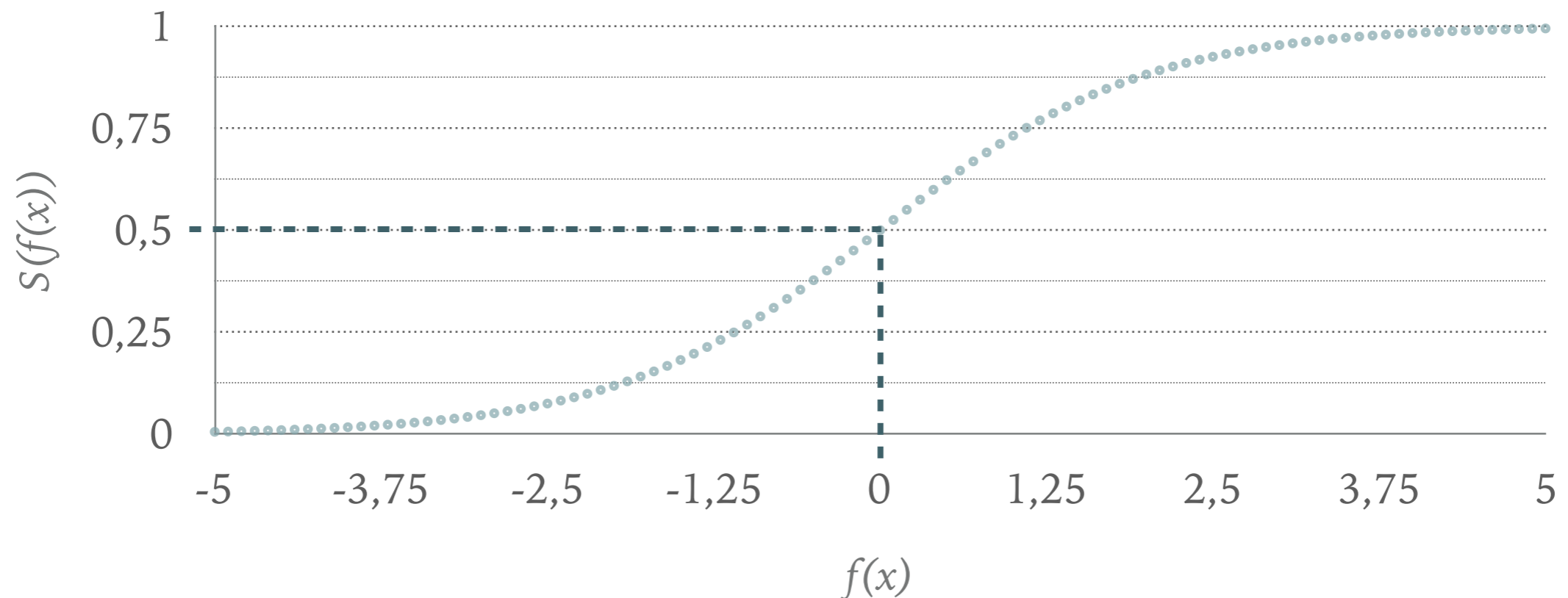
MODEL

- Decision rule

- $\hat{y} = \operatorname{argmax}_{y \in \{0,1\}} (P(Y = y | X = z))$

- Linear classifier

- Hyperplane equation $f(x) = 0$



PARAMETER ESTIMATION

➤ Maximum likelihood estimation

$$\begin{aligned} \text{➤ } L(D; w) &= \prod_{i=1}^n P(Y = y_i | X = x_i) \\ &= \prod_{i=1}^n p(1 | x_i)^{y_i} (1 - p(1 | x_i))^{1-y_i} \end{aligned}$$

$$\begin{aligned} \text{➤ } \ell(D; w) &= \log \left(\prod_{i=1}^n p(1 | x_i)^{y_i} (1 - p(1 | x_i))^{1-y_i} \right) \\ &= \sum_{i=1}^n y_i \log(p(1 | x_i)) + (1 - y_i) \log((1 - p(1 | x_i))) \end{aligned}$$

➤ Problem formulation

$$\text{➤ } w^* = \underset{w}{\operatorname{argmax}}(\ell(D; w)) = \underset{w}{\operatorname{argmax}}(L(D; w))$$

PARAMETER ESTIMATION

- Rewrite the model and incorporate a constant $x_0 = 1$ so that

- $$P(Y = 1 | X = x) = \frac{1}{1 + e^{-wx}}$$
$$= S(wx)$$

- Partial derivative of the log-likelihood in w_j

- $$\frac{\partial}{\partial w_j} \ell(D; w) = \frac{\partial}{\partial w_j} \sum_{i=1}^n y_i \log(p(1 | x_i)) + (1 - y_i) \log((1 - p(1 | x_i)))$$
$$= \frac{\partial}{\partial w_j} \sum_{i=1}^n y_i \log(S(wx_i)) + (1 - y_i) \log((1 - S(wx_i)))$$
$$= \sum_{i=1}^n y_i \frac{\partial}{\partial w_j} \log(S(wx_i)) + (1 - y_i) \frac{\partial}{\partial w_j} \log((1 - S(wx_i)))$$

PARAMETER ESTIMATION

- ▶ Partial derivative of the left-hand side of the log-likelihood

$$\begin{aligned} \text{▶ } \frac{\partial}{\partial w_j} \log(S(wx_i)) &= \frac{\partial}{\partial w_j} \log\left(\frac{1}{1 + e^{-wx_i}}\right) \\ &= \frac{\partial}{\partial w_j} \log(1) - \log(1 + e^{-wx_i}) \\ &= \frac{\partial}{\partial w_j} - \log(1 + e^{-wx_i}) \\ &= - \frac{-x_{ij}e^{-wx_i}}{1 + e^{-wx_i}} \\ &= x_{ij} \left(1 - \frac{1}{1 + e^{-wx_i}}\right) \\ &= x_{ij}(1 - S(wx_i)) \end{aligned}$$

PARAMETER ESTIMATION

► Partial derivative of the right-hand side of the log-likelihood

$$\begin{aligned} \text{► } \frac{\partial}{\partial w_j} \log(1 - S(wx_i)) &= \frac{\partial}{\partial w_j} \log\left(1 - \frac{1}{1 + e^{-wx_i}}\right) \\ &= \frac{\partial}{\partial w_j} \log\left(\frac{1 + e^{-wx_i}}{1 + e^{-wx_i}} - \frac{1}{1 + e^{-wx_i}}\right) \\ &= \frac{\partial}{\partial w_j} \log\left(\frac{e^{-wx_i}}{1 + e^{-wx_i}}\right) \\ &= \frac{\partial}{\partial w_j} \log(e^{-wx_i}) - \log(1 + e^{-wx_i}) \\ &= \frac{\partial}{\partial w_j} -wx_i - \log(1 + e^{-wx_i}) \\ &= -x_{ij} + x_{ij}(1 - S(wx_i)) \\ &= -x_{ij}S(wx_i) \end{aligned}$$

PARAMETER ESTIMATION

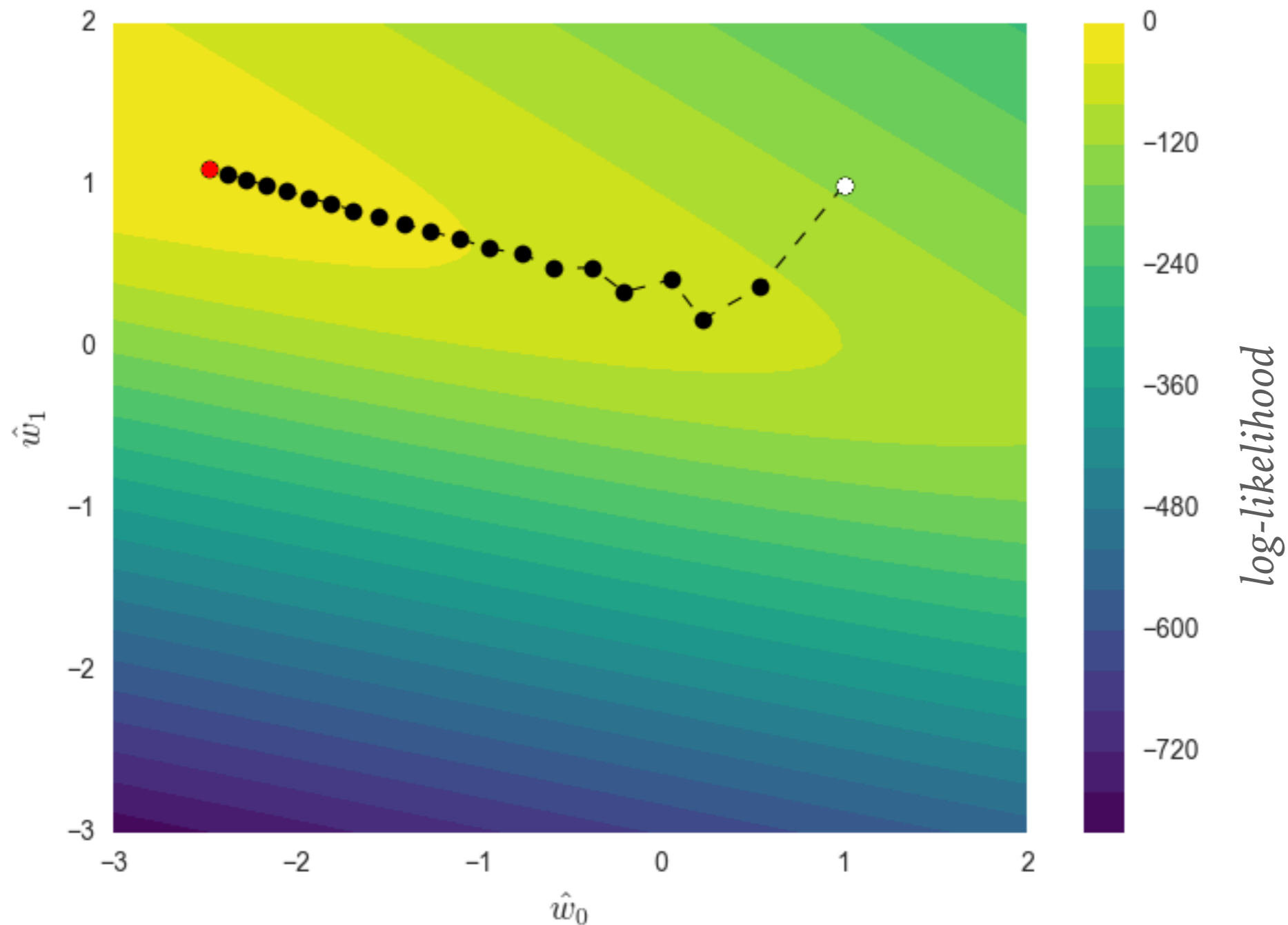
- Assemble the two partial derivatives

$$\begin{aligned}\text{➤ } \frac{\partial}{\partial w_j} \ell &= \sum_{i=1}^n y_i x_{ij} (1 - S(wx_i)) - x_{ij} (1 - y_i) S(wx_i) \\ &= \sum_{i=1}^n y_i x_{ij} - y_i x_{ij} S(wx_i) - x_{ij} S(wx_i) + y_i x_{ij} S(wx_i) \\ &= \sum_{i=1}^n y_i x_{ij} - x_{ij} S(wx_i) \\ &= \sum_{i=1}^n x_{ij} (y_i - S(wx_i))\end{aligned}$$

- There's no analytic solution to $\frac{\partial}{\partial w_j} \ell(D; w) = 0$

APPROXIMATION VIA GRADIENT DESCENT

- Simple example of gradient descent



REGULARIZATION

- Add a penalization term, r , to the function to maximize

- $w^* = \underset{w}{\operatorname{argmax}} (\ell(D; w) - \lambda r(w))$

- Common penalization terms

- Ridge

- $w^* = \underset{w}{\operatorname{argmax}} (\ell(D; w) - \lambda \|w\|_2^2)$ with $\|w\|_2^2 = \sum_{j=0}^d w_j^2$

- Penalize large coefficients

- LASSO

- $w^* = \underset{w}{\operatorname{argmax}} (\ell(D; w) - \lambda \|w\|_1)$ with $\|w\|_1 = \sum_{j=0}^d |w_j|$

- Favor sparse solutions

EVALUATION

CROSS-VALIDATION

- Split the data into two sets
 - Training set
 - Test set
- Methodology
 - Estimate the model parameters on the train set
 - Predict on the test set
 - Build the confusion matrix, M , *i.e.* a contingency table to compare predicted labels and true labels
- Repeat this k times with different training and test sets
 - k -fold cross-validation

CROSS-VALIDATION

- Common evaluation metrics

- Overall accuracy

- $A = \frac{\text{\# of correctly classified documents}}{\text{total \# of documents}}$

- Precision w.r.t a given class

- $P_c = \frac{\text{\# of documents correctly classified in } c}{\text{\# of documents classified in } c}$

- Recall w.r.t a given class

- $R_c = \frac{\text{\# of documents correctly classified in } c}{\text{\# of documents belonging to class } c}$

TAKE AWAY MESSAGE

- The multinomial naïve Bayes classifier is a good baseline, always try it first
 - Efficient parameter estimation
 - Sensitive to vocabulary filtering
 - Limited to multinomial data, *i.e.* raw word counts
- The classifier based on the logistic regression is more versatile
 - Compatible with any weighting scheme
 - Useful regularization
 - Higher computational complexity
- Both are linear classifiers