

Machine Learning

Lecture 8: Linear Classification

Prof. Dr. Aleksandar Bojchevski

14.05.24

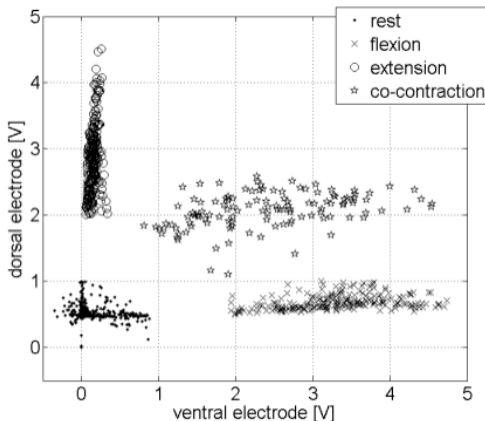
Hard deterministic classification

Probabilistic generative models for linear classification

Probabilistic discriminative models for linear classification

Classification problem

Given features $\mathbf{x}_i \in \mathbb{R}^D$ and labels $y_i \in \mathcal{C} = \{1, \dots, C\}$, find a mapping $f : \mathbb{R}^D \rightarrow \mathcal{C}$ such that $y_i = f(\mathbf{x}_i)$ for as many examples as possible.



Zero-one loss

How do we measure the quality of a prediction $\hat{y} = f(\mathbf{x})$?

The zero-one loss $l_{01}(y, \hat{y}) = \mathbb{I}(y \neq \hat{y})$ equals 1 if we have a misclassification.

Given a dataset $\mathcal{D} = \{(\mathbf{x}, y_i)\}_{i=1}^N$, the total number of misclassified examples is

$$\mathcal{L} = \sum_{i=1}^N l_{01}(y, \hat{y}) = \sum_{i=1}^N \mathbb{I}(\hat{y}_i \neq y_i) = \mathbb{I}(f(\mathbf{x})_i \neq y_i)$$

How do we choose a good $f(\cdot)$?

Hyperplane as a decision boundary

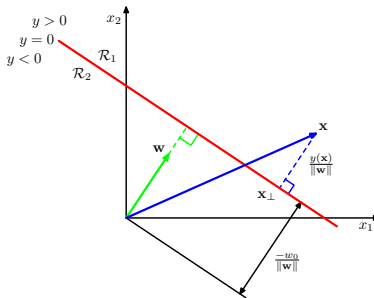
For a 2-class problem, $\mathcal{C} = \{0, 1\}$, we can try to separate the instances from the two classes by a [hyperplane](#).

Hyperplane as a decision boundary

For a 2-class problem, $\mathcal{C} = \{0, 1\}$, we can try to separate the instances from the two classes by a [hyperplane](#).

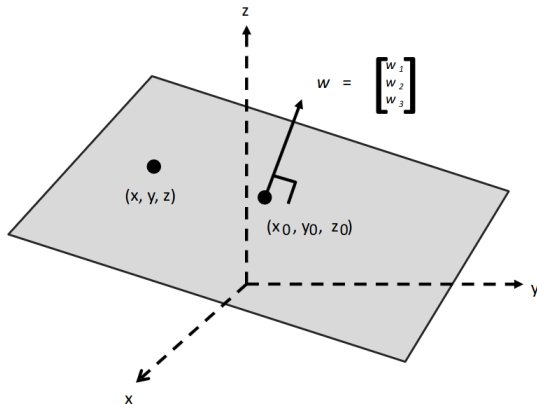
A hyperplane be defined by a normal vector \mathbf{w} and an offset w_0 .
Hyperplanes are computationally very convenient: easy to evaluate.

$$\mathbf{w}^T \mathbf{x} + w_0 \begin{cases} = 0 & \text{if } \mathbf{x} \text{ on the plane} \\ > 0 & \text{if } \mathbf{x} \text{ on normal's side} \\ < 0 & \text{else} \end{cases}$$



Hyperplane as a decision boundary

Same idea applies in higher dimensions. We have $\mathbf{w}^T(\mathbf{x} - \mathbf{x}_0) = \mathbf{w}^T \mathbf{x} - \underbrace{\mathbf{w}^T \mathbf{x}_0}_{=w_0} = 0$.



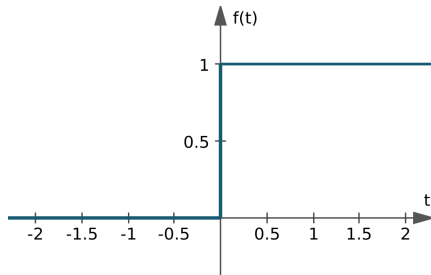
Making predictions

Decision rule

$$\hat{y} = f(\mathbf{w}^T \mathbf{x} + w_0)$$

where f is the step function defined as:

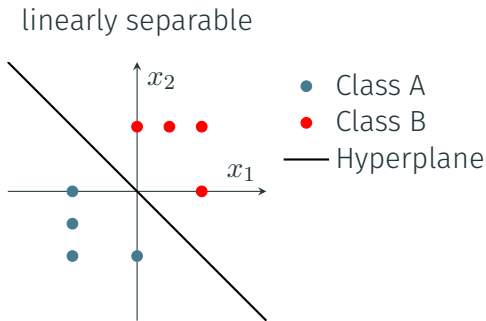
$$f(t) = \begin{cases} 1 & \text{if } t > 0, \\ 0 & \text{otherwise.} \end{cases}$$



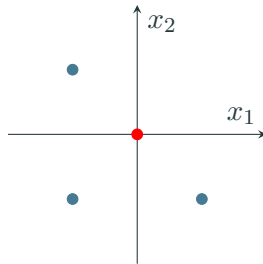
Sometimes the labels are equivalently given as $y_i \in \{1, -1\}$ and the predictions are $\hat{y}_i = \text{sign}(f(\mathbf{x}_i))$ where $\text{sign}(\cdot) = 1$ if the input is positive and -1 otherwise.

Linear separability

A data set $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$ is **linearly separable** if there exists a hyperplane for which all \mathbf{x}_i with $y_i = 0$ are on one side, and all \mathbf{x}_i with $y_i = 1$ on the other side.



not linearly separable



Learning rule for the perceptron

The perceptron algorithm is one of the oldest methods for binary classification.

Initialize parameters to any value, e.g., a zero vector: $\boldsymbol{w}, w_0 \leftarrow \mathbf{0}$.

¹However, there is no way to determine the number of required iterations in advance.

Learning rule for the perceptron

The perceptron algorithm is one of the oldest methods for binary classification.

Initialize parameters to any value, e.g., a zero vector: $\mathbf{w}, w_0 \leftarrow \mathbf{0}$.

For each misclassified sample \mathbf{x}_i in the training set update

$$\mathbf{w} \leftarrow \begin{cases} \mathbf{w} + \mathbf{x}_i & \text{if } y_i = 1, \\ \mathbf{w} - \mathbf{x}_i & \text{if } y_i = 0. \end{cases} \quad w_0 \leftarrow \begin{cases} w_0 + 1 & \text{if } y_i = 1, \\ w_0 - 1 & \text{if } y_i = 0. \end{cases}$$

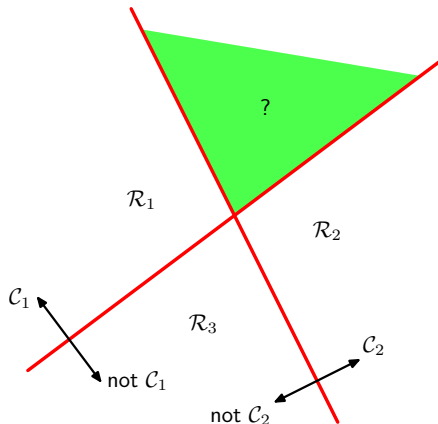
until all samples are classified correctly.

This method takes a finite¹ number of steps to converge to a (\mathbf{w}, w_0) discriminating between two classes **if it exists** (if the data is linearly separable).

¹However, there is no way to determine the number of required iterations in advance.

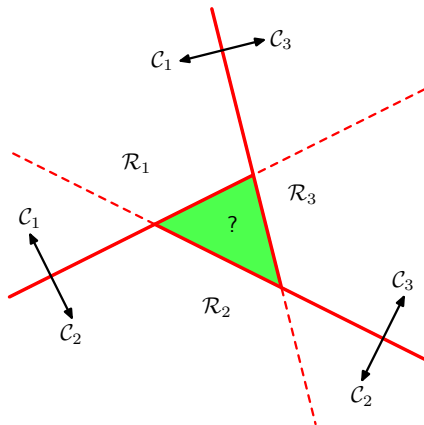
Handling multiple classes

One-versus-rest classifier: Each hyperplane \mathcal{H}_i decides class \mathcal{C}_i vs. not class \mathcal{C}_i .



Handling multiple classes

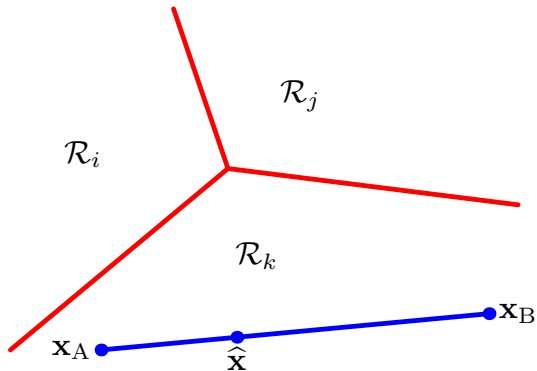
One-versus-one classifier: hyperplane \mathcal{H}_{ij} decides each pair class \mathcal{C}_i vs. class \mathcal{C}_j .
Use majority vote to classify.



Handling multiple classes

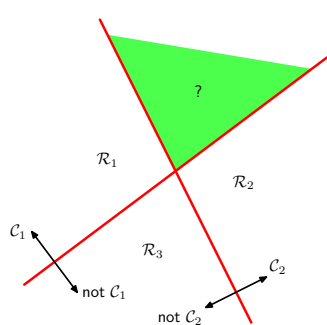
Multiclass discriminant: define C linear functions $f_c(\mathbf{x}) = \mathbf{w}_c^T \mathbf{x} + w_{0c}$ with the decision rule $\hat{y} = \arg \max_{c \in \mathcal{C}} f_c(\mathbf{x})$.

That is, assign \mathbf{x} to the class c which produces the highest $f_c(\mathbf{x})$, dividing the domain into convex decision regions \mathcal{R}_i .

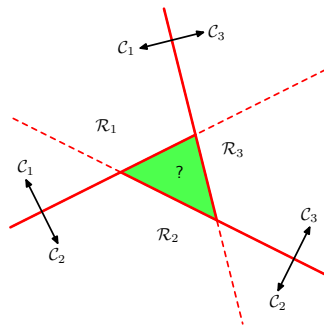


Handling multiple classes

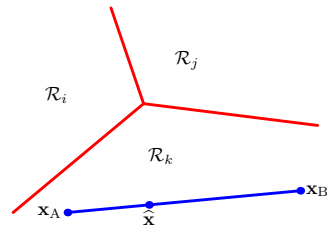
Usually we prefer the multiclass discriminant unless we have no choice.



One-vs-one

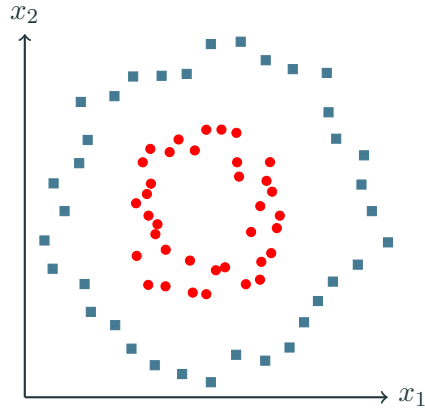


One-vs-rest



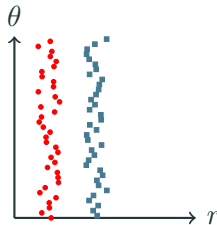
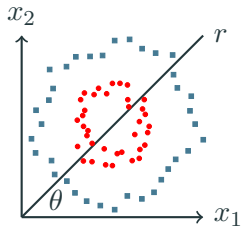
Multiclass

What if the classes are not linearly separable?



Basis functions

Apply a nonlinear transformation $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$ that maps samples to a space where they are linearly separable.

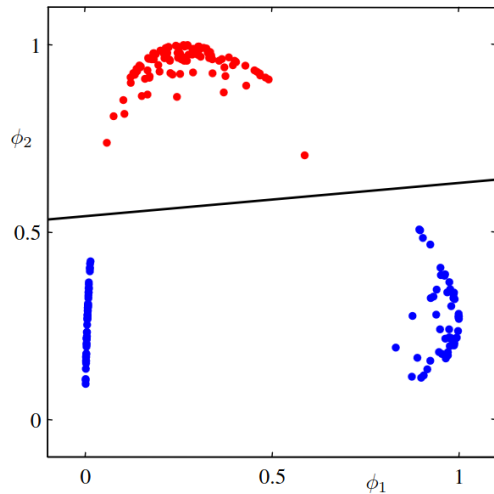
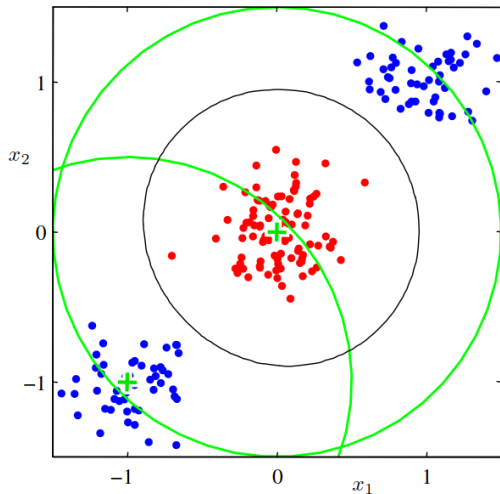


Here, $\phi(\mathbf{x}) = (r, \theta) = (\|\mathbf{x}\|_2, \text{angle}(\mathbf{x}))$.

Here we talk about fixed basis functions. A deeper discussion of this method will take place in *Kernels*. Adaptive basis functions will be covered in *Deep Learning*.

Basis functions

Linear separability using two Gaussian basis functions $\phi_1(\mathbf{x})$, $\phi_2(\mathbf{x})$ (green).



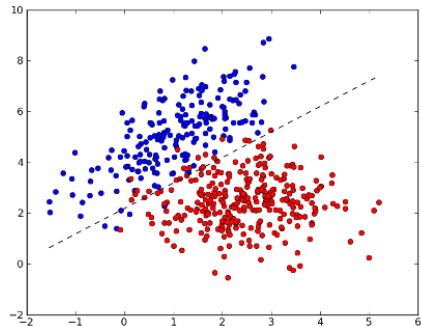
Limitations of hard-decision based classifiers

No measure of uncertainty

Cannot handle noisy data

Poor generalization

Difficult to optimize



What are the alternatives?

Probabilistic models for classification

Solution: model the distribution of the class label y given the data \mathbf{x} .

$$p(y = c \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid y = c) \cdot p(y = c)}{p(\mathbf{x})}$$

Two types of models:

Generative

- Model the joint distribution $p(\mathbf{x}, y = c) = p(\mathbf{x} \mid y = c) \cdot p(y = c)$

Discriminative

- Directly model the distribution $p(y = c \mid \mathbf{x})$

Given $p(y \mid \mathbf{x})$ we can make the prediction \hat{y} based on our problem.

Popular choice is the mode: $\hat{y} = \arg \max_{c \in \mathcal{C}} p(y = c \mid \mathbf{x})$.

Hard deterministic classification

Probabilistic generative models for linear classification

Probabilistic discriminative models for linear classification

The idea is to obtain the class posterior using Bayes' theorem

$$p(y = c \mid \mathbf{x}) \propto \underbrace{p(\mathbf{x} \mid y = c)}_{\text{class conditional}} \cdot \underbrace{p(y = c)}_{\text{class prior}}$$

The model consists of:

class prior - a priori probability of a point belonging to a class c ,

class conditional - probability of observing \mathbf{x} , given that it belongs to class c .

Applying a generative model

Choose a parametric model for the class conditional $p(\mathbf{x} \mid y = c, \boldsymbol{\psi})$ and the class prior $p(y = c \mid \boldsymbol{\theta})$.

Learning equals estimating the parameters of our model $\{\boldsymbol{\psi}, \boldsymbol{\theta}\}$ from the data \mathcal{D} , e.g., using maximum likelihood to obtain estimates $\{\hat{\boldsymbol{\psi}}, \hat{\boldsymbol{\theta}}\}$.

Once fitted, we can perform **inference** - classify a new \mathbf{x}_{new} using Bayes rule

$$p(y_{\text{new}} = c \mid \mathbf{x}_{\text{new}}, \hat{\boldsymbol{\psi}}, \hat{\boldsymbol{\theta}}) \propto p(\mathbf{x}_{\text{new}} \mid y_{\text{new}} = c, \hat{\boldsymbol{\psi}})p(y_{\text{new}} = c \mid \hat{\boldsymbol{\theta}})$$

Additionally, we can **generate** new data - hence the name:

- sample a class label $y_{\text{sample}} \sim p(y \mid \hat{\boldsymbol{\theta}})$
- sample a feature vector $\mathbf{x}_{\text{sample}} \sim p(\mathbf{x} \mid y = y_{\text{sample}}, \hat{\boldsymbol{\psi}})$

How do we choose the class prior $p(y = c)$?

The label y can take one of C discrete values \implies use categorical distribution.

$$y \sim \text{Categorical}(\boldsymbol{\theta})$$

The parameter $\boldsymbol{\theta} \in \mathbb{R}^C$ specifies the probability of each class

$$p(y = c) = \theta_c \quad \text{or equivalently} \quad p(y) = \prod_{c=1}^C \theta_c^{\mathbb{I}(y=c)}$$

and is subject to the constraints $0 \leq \theta_c \leq 1$ and $\sum_{c=1}^C \theta_c = 1$.

The maximum likelihood estimate $\hat{\boldsymbol{\theta}}$ for $\boldsymbol{\theta}$ given the data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$ is

How do we choose the class prior $p(y = c)$?

The label y can take one of C discrete values \implies use categorical distribution.

$$y \sim \text{Categorical}(\boldsymbol{\theta})$$

The parameter $\boldsymbol{\theta} \in \mathbb{R}^C$ specifies the probability of each class

$$p(y = c) = \theta_c \quad \text{or equivalently} \quad p(y) = \prod_{c=1}^C \theta_c^{\mathbb{I}(y=c)}$$

and is subject to the constraints $0 \leq \theta_c \leq 1$ and $\sum_{c=1}^C \theta_c = 1$.

The maximum likelihood estimate $\hat{\boldsymbol{\theta}}$ for $\boldsymbol{\theta}$ given the data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$ is

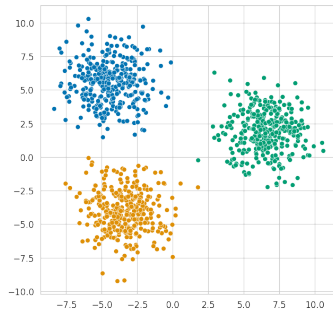
$$\hat{\theta}_c = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i = c)$$

How do we choose the class conditionals $p(\mathbf{x} \mid y = c)$?

For continuous features, $\mathbf{x} \in \mathbb{R}^D$, we can use a multivariate normal for each class.

$$p(\mathbf{x} \mid y = c) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_c, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) \right\}$$

We use the same $\boldsymbol{\Sigma}$ for each class, as estimating all $\boldsymbol{\Sigma}_c$'s behaves badly numerically, unless we have **lots** of data.



We will derive MLE estimates for $\{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_C, \boldsymbol{\Sigma}\}$ in the tutorial (see Bishop 4.2.2).

Posterior distribution

Now that we have chosen $p(\mathbf{x} \mid y)$ and $p(y)$, and have estimated their parameters from the training data, how do we perform classification?

Let's assume for simplicity that we have two classes $\mathcal{C} = \{0, 1\}$.

$$\begin{aligned} p(y = 1 \mid \mathbf{x}) &= \frac{p(\mathbf{x} \mid y = 1) p(y = 1)}{p(\mathbf{x} \mid y = 1) p(y = 1) + p(\mathbf{x} \mid y = 0) p(y = 0)} \\ &= \frac{1}{1 + \exp(-a)} =: \sigma(a) \end{aligned}$$

where σ is the **sigmoid** function and we defined

$$a = \log \frac{p(\mathbf{x} \mid y = 1) p(y = 1)}{p(\mathbf{x} \mid y = 0) p(y = 0)}$$

To avoid clutter, we implicitly condition the distributions on their respective parameters $(\boldsymbol{\theta}, \boldsymbol{\mu}_c, \boldsymbol{\Sigma})$.

Linear discriminant analysis (LDA)

For Gaussian class-conditionals with the same covariance Σ we have

$$\begin{aligned}a &= \log \frac{p(\mathbf{x} \mid y = 1)p(y = 1)}{p(\mathbf{x} \mid y = 0)p(y = 0)} \\&= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) + \log p(y = 1) \\&\quad + \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_0)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_0) - \log p(y = 0) \\&= \mathbf{w}^T \mathbf{x} + w_0\end{aligned}$$

where we define

$$\begin{aligned}\mathbf{w} &= \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) \\w_0 &= -\frac{1}{2}\boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_0^T \Sigma^{-1} \boldsymbol{\mu}_0 + \log \frac{p(y = 1)}{p(y = 0)}\end{aligned}$$

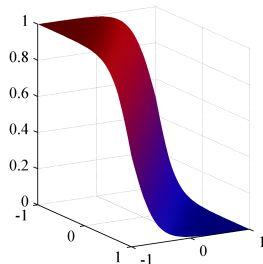
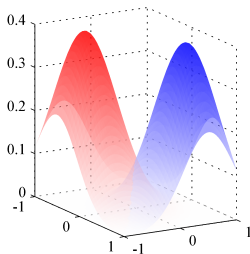
LDA for $C = 2$ classes

This means, that the posterior distribution is a sigmoid of a linear function of \mathbf{x}

$$p(y = 1 \mid \mathbf{x}) = \frac{1}{1 + \exp(-(\mathbf{w}^T \mathbf{x} + w_0))} = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

Or equivalently $y \mid \mathbf{x} \sim \text{Bernoulli}(\sigma(\mathbf{w}^T \mathbf{x} + w_0))$.

For $D = 2$, $p(\mathbf{x} \mid y = 1)$ in red, and $p(\mathbf{x} \mid y = 0)$ in blue. Right shows $p(y = 1 \mid \mathbf{x})$.



LDA for $C > 2$ classes

Using Bayes' theorem, the posterior for the $C > 2$ case is

$$p(y = c \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid y = c)p(y = c)}{\sum_{c'=1}^C p(\mathbf{x} \mid y = c')p(y = c')}.$$

Working out the math, we get

$$= \frac{\exp(\mathbf{w}_c^T \mathbf{x} + w_{c0})}{\sum_{c'=1}^C \exp(\mathbf{w}_{c'}^T \mathbf{x} + w_{c'0})}$$

where

$$\begin{aligned}\mathbf{w}_c &= \Sigma^{-1} \boldsymbol{\mu}_c \\ w_{c0} &= -\frac{1}{2} \boldsymbol{\mu}_c^T \Sigma^{-1} \boldsymbol{\mu}_c + \log p(y = c)\end{aligned}$$

Softmax function

On the previous slide we arrived at the `softmax` function defined as

$$\sigma(\mathbf{x})_i = \frac{\exp(x_i)}{\sum_{k=1}^K \exp(x_k)}$$

Softmax σ is a generalization of sigmoid to multiple dimensions

$$\sigma : \mathbb{R}^K \rightarrow \Delta^{K-1}$$

where

$$\Delta^{K-1} = \left\{ \mathbf{x} \in \mathbb{R}^K \mid \sum_{k=1}^K x_k = 1 \text{ and } x_k \geq 0, k = 1, \dots, K \right\}$$

is the standard **probability simplex**.

Class conditionals: Variant II (Naive Bayes)

Naive Bayes: Assume that the D features of a sample $\mathbf{x} = (x_1, x_2, \dots, x_d)$ are conditionally independent given the class, i.e.

$$p(x_1, x_2, \dots, x_d | y = c) = \prod_{i=1}^d p(x_i | y = c)$$

In the case of continuous data where the likelihood is assumed to be a normal distribution, this corresponds to **diagonal** covariance matrices, i.e.

$$p(\mathbf{x} | y = c) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$$

Important: We use a different (diagonal) covariance matrix $\boldsymbol{\Sigma}_c$ for each class c . Compare to LDA where all class conditionals share the same covariance matrix $\boldsymbol{\Sigma}$.

Assume two classes $\mathcal{C} = \{0, 1\}$. Like in LDA we need to compute

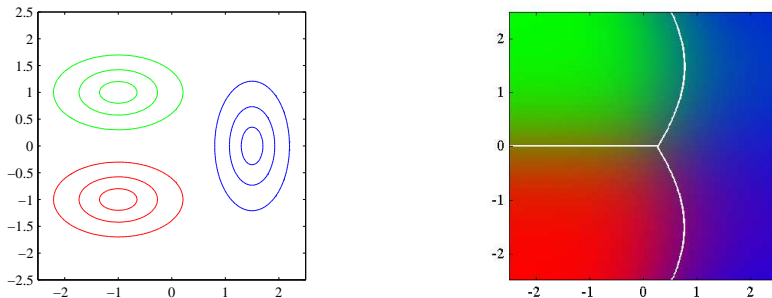
$$\begin{aligned}a &= \log \frac{p(\mathbf{x} \mid y = 1)p(y = 1)}{p(\mathbf{x} \mid y = 0)p(y = 0)} \\&= \frac{1}{2}\mathbf{x}^T[\boldsymbol{\Sigma}_0^{-1} - \boldsymbol{\Sigma}_1^{-1}]\mathbf{x} + \mathbf{x}^T[\boldsymbol{\Sigma}_1^{-1}\boldsymbol{\mu}_1 - \boldsymbol{\Sigma}_0^{-1}\boldsymbol{\mu}_0] \\&\quad - \frac{1}{2}\boldsymbol{\mu}_1^T\boldsymbol{\Sigma}_1^{-1}\boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_0^T\boldsymbol{\Sigma}_0^{-1}\boldsymbol{\mu}_0 + \log \frac{\pi_1}{\pi_0} + \frac{1}{2} \log \frac{|\boldsymbol{\Sigma}_0|}{|\boldsymbol{\Sigma}_1|} \\&= \mathbf{x}^T\mathbf{W}_2\mathbf{x} + \mathbf{w}_1^T\mathbf{x} + w_0\end{aligned}$$

where we define

$$\begin{aligned}\mathbf{W}_2 &= \frac{1}{2}[\boldsymbol{\Sigma}_0^{-1} - \boldsymbol{\Sigma}_1^{-1}] \\ \mathbf{w}_1 &= \boldsymbol{\Sigma}_1^{-1}\boldsymbol{\mu}_1 - \boldsymbol{\Sigma}_0^{-1}\boldsymbol{\mu}_0 \\ w_0 &= -\frac{1}{2}\boldsymbol{\mu}_1^T\boldsymbol{\Sigma}_1^{-1}\boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_0^T\boldsymbol{\Sigma}_0^{-1}\boldsymbol{\mu}_0 + \log \frac{\pi_1}{\pi_0} + \frac{1}{2} \log \frac{|\boldsymbol{\Sigma}_0|}{|\boldsymbol{\Sigma}_1|}\end{aligned}$$

Naive Bayes: Decision Boundary

Left: Gaussian distributions with diagonal covariance, right: decision boundary.



Naive Bayes results in a quadratic decision boundary.
Compare to LDA which leads to a linear decision boundary.

Naive Bayes: Advantage

Class conditional for Naive Bayes: $p(x_1, x_2, \dots, x_d | y = c) = \prod_{i=1}^d p(x_i | y = c)$.

Advantage: feature independence allows us to easily handle different data types.

Simply choose a suitable (univariate) distribution for each feature:

x_1 – Gaussian, x_2 – Categorical, ..., etc.

Hard deterministic classification

Probabilistic generative models for linear classification

Probabilistic discriminative models for linear classification

Probabilistic discriminative model

An alternative approach to generative modeling is to model the posterior distribution $p(y \mid \mathbf{x})$ directly. Such models are called **discriminative**.

We saw that a generative approach with Gaussian class-conditionals and shared covariance matrix Σ (LDA) leads to the posterior distribution

$$\begin{aligned} p(y = 1 \mid \mathbf{x}) &= \sigma(\mathbf{x}^T \mathbf{w} + w_0), \\ p(y = 0 \mid \mathbf{x}) &= 1 - \sigma(\mathbf{x}^T \mathbf{w} + w_0) \end{aligned}$$

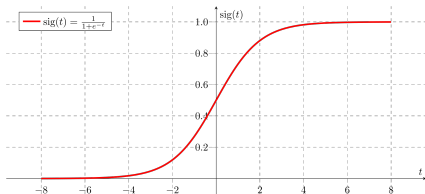
where \mathbf{w}, w_0 depend on the parameters of class-conditionals μ_0, μ_1, Σ .

Why not just let \mathbf{w} and w_0 be free parameters and choose them directly?

Logistic regression

We model the posterior distribution as $y \mid \mathbf{x} \sim \text{Bernoulli}(\sigma(\mathbf{w}^T \mathbf{x} + w_0))$

where $\sigma(a) = \frac{1}{1+\exp(-a)}$ is the sigmoid and \mathbf{w}, w_0 are free model parameters.



This model is called **logistic regression**.

We again absorb the bias term by overloading the notation and defining $\mathbf{w}^T \mathbf{x} := w_0 + w_1 x_1 + \dots + w_D x_D$ which is equivalent to defining $x_0 = 1$.

Likelihood of the logistic regression

Assuming that all samples (\mathbf{x}_i, y_i) are drawn i.i.d., we can write the likelihood as

$$\begin{aligned} p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}) &= \prod_{i=1}^N p(y_i \mid \mathbf{x}_i, \mathbf{w}) \\ &= \prod_{i=1}^N \underbrace{p(y = 1 \mid \mathbf{x}_i, \mathbf{w})^{y_i}}_{=1 \text{ if } y_i=0} \underbrace{(1 - p(y = 1 \mid \mathbf{x}_i, \mathbf{w}))^{1-y_i}}_{=1 \text{ if } y_i=1} \\ &= \prod_{i=1}^N \sigma(\mathbf{w}^T \mathbf{x}_i)^{y_i} (1 - \sigma(\mathbf{w}^T \mathbf{x}_i))^{1-y_i} \end{aligned}$$

Negative log-likelihood

Similar to linear regression, we can define the **negative log-likelihood** as the loss:

$$\begin{aligned}\mathcal{L}(\mathbf{w}) &= -\log p(\mathbf{y} \mid \mathbf{w}, \mathbf{X}) \\ &= -\sum_{i=1}^N (y_i \log \sigma(\mathbf{x}_i) + (1 - y_i) \log(1 - \sigma(\mathbf{x}_i)))\end{aligned}$$

This loss function is also called the **binary cross entropy**.

Finding the maximum likelihood estimate for \mathbf{w} is equivalent to solving

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w})$$

Solving the minimization problem

There is **no closed-form** solution for logistic regression.

This means, we cannot represent the optimal \mathbf{w}^* directly using standard mathematical operations, such as multiplication, matrix inversion, etc.

However, there is still hope! We can use standard optimization techniques to numerically solve our problem. We will cover this in the upcoming lectures.

For now, just assume that we can find \mathbf{w}^* .

Logistic regression + weights regularization

Since maximum likelihood estimation can lead to overfitting, just like in linear regression, we can control it by penalizing large weights.

$$\mathcal{L}_{\text{reg}}(\mathbf{w}) = -\log p(\mathbf{y} \mid \mathbf{w}, \mathbf{X}) + \lambda \|\mathbf{w}\|_q^q$$

As before, $q = 2$ corresponds to MAP estimation with a Gaussian prior on \mathbf{w} .

Again, there is no closed-form solution available.

Multiclass logistic regression

For the binary classification we used the sigmoid function to “squeeze” the unnormalized probability $\mathbf{w}^T \mathbf{x}$ into the range $(0, 1)$.

The same can be done for multiple classes using the [softmax](#) function.

$$p(y = c \mid \mathbf{x}) = \frac{\exp(\mathbf{w}_c^T \mathbf{x})}{\sum_{c'} \exp(\mathbf{w}_{c'}^T \mathbf{x})}$$

How does this relate to multiclass LDA?

Loss for multiclass logistic regression

The negative log-likelihood for multiclass logistic regression can be written as

$$\begin{aligned}\mathcal{L}(\mathbf{w}) &= -\log p(\mathbf{Y} \mid \mathbf{X}, \mathbf{w}) \\ &= -\sum_{i=1}^N \sum_{c=1}^C y_{ic} \log p(y_i = c \mid \mathbf{x}_i, \mathbf{w}) \\ &= -\sum_{i=1}^N \sum_{c=1}^C y_{ic} \log \frac{\exp(\mathbf{w}_c^T \mathbf{x})}{\sum_{c'} \exp(\mathbf{w}_{c'}^T \mathbf{x})}\end{aligned}$$

and is also called the **cross entropy**.

Here we use **one-hot encoding**: vector of categorical variables $\mathbf{y} \in \mathcal{C}^N$ is encoded as a binary matrix $\mathbf{Y} \in \{0, 1\}^{N \times C}$, where $y_{ic} = \begin{cases} 1 & \text{if sample } i \text{ belongs to class } c \\ 0 & \text{else} \end{cases}$

Generative vs. discriminative models

Discriminative models usually tend to achieve better performance² when it comes to pure classification tasks.

While generative models work reasonably well when their assumptions hold, they can be fragile when these assumptions are violated.

Nevertheless, generative models provide the added benefits of better handling missing data, detecting outliers, generating new data and being more appropriate in the semi-supervised setting.

²If we measure performance just by the accuracy, however, this is not necessarily true for other metrics such as the calibration error or fairness.

Main reading

- "Pattern Recognition and Machine Learning" by Bishop
[ch. 4.1.1, 4.1.2, 4.1.7, 4.2, 4.3.0–4.3.4]

Slides are based on an older version by G. Jensen, C. Osendorfer, and S. Günnemann. Some figures are from Bishop's "Pattern Recognition and Machine Learning".