Linear Classifiers

Volker Tresp Winter 2023-2024

Classification

- Classification is the central task of pattern recognition
- Sensors supply information about an object: to which class does the object belong (dog, cat, ...)?

Overlapping Classes

- The beauty of Machine Learning is that a few model classes (neural networks, kernel approaches, ...) can be applied to almost any supervised learning task
- This hides a bit that the data settings can be quite different
- There are problems where class boundaries are well defined but maybe quite complex; an example is OCR; here Deep Neural Networks, manifold learning and kernel systems are quite effective; this concerns often our Cases I and II
- In other applications there is little structure in the data and classes overlap; this the situation encountered in many healthcare applications (biomedicine); this concerns often our Cases III and IV
- Often, the problem is not as much to separate classes, but to show that there is a signal at all; the question might be if there is a detectable positive effect of the new medication!

Linear Classifiers

- Linear classifiers separate classes by a linear hyperplane
- In high dimensions a linear classifier often can separate the classes
- Linear classifiers cannot solve the *exclusive-or* problem
- In combination with basis functions, kernels or a neural network, linear classifiers can form nonlinear class boundaries

Hard and Soft (sigmoid) Transfer Functions



• First, the activation function of the neurons in the hidden layer are calculated as the weighted sum of the inputs x_i as

$$h(\mathbf{x}) = \sum_{j=0}^{M} w_j x_j$$

(note: $x_0 = 1$ is a constant input, so that w_0 corresponds to the bias)

- The sigmoid neuron has a soft (sigmoid) transfer function
 - Perceptron : $\hat{y} = sign(h(\mathbf{x}))$

Sigmoid function: $\hat{y} = sig(h(\mathbf{x}))$

Binary Classification Problems

- We will focus first on binary classification where the task is to assign binary class labels $y_i = 1$ and $y_i = 0$ (or $y_i = 1$ and $y_i = -1$)
- We already know the *Perceptron*. Now we learn about additional approaches
 - I. Generative models for classification
 - II. Logistic regression
 - III. Classification via regression

Two Linearly Separable Classes



Two Classes that Cannot be Separated Linearly



The Classical Example not two Classes that cannot be Separated Linearly: XOR



Separability is not a Goal in Itself. With Overlapping Classes the Goal is the Best Possible Hyperplane



I. Generative Model for Classification

- In a generative model one assumes a probabilistic data generating process (likelihood model). Often generative models are complex and contain unobserved (latent, hidden) variables
- Here we consider a simple example: data is generated from class-specific Gaussian distributions
- First we have a model how classes are generated P(y). y = 1 could stand for a good customer and y = 0 could stand for a bad customer.

Generative Model for Classification (cont'd)

- Then we have a model how attributes are generated, given the classes P(x
 y). This could stand for
 - Income, age, occupation $(\tilde{\mathbf{x}})$ given a customer is a good customer (y = 1)
 - Income, age, occupation $(\tilde{\mathbf{x}})$ given a customer is not a good customer (y = 0)
- Using Bayes formula, we then derive P(y|x): the probability that a given customer is
 a good customer y = 1 or bad customer y = 0, given that we know the customer's
 income, age and occupation

How is Data Generated?

• We assume that the observed classes y_i are generated with probability

$$P(y_i = 1) = \kappa_1$$
 $P(y_i = 0) = \kappa_0 = 1 - \kappa_1$

with $0 \leq \kappa_1 \leq 1$.

- In a next step, a data point $ilde{\mathbf{x}}_i$ has been generated from $P(ilde{\mathbf{x}}_i|y_i)$
- (Note, that $\tilde{\mathbf{x}}_i = (x_{i,1}, \dots, x_{i,M})^T$, which means that $\tilde{\mathbf{x}}_i$ does not contain the bias $x_{i,0}$)
- We now have a complete model: $P(y_i)P(ilde{\mathbf{x}}_i|y_i)$

Bayes' Theorem

• To classify a data point $\tilde{\mathbf{x}}_i$, i.e. to determine the y_i , we apply Bayes theorem and get

$$P(y_i|\tilde{\mathbf{x}}_i) = \frac{P(\tilde{\mathbf{x}}_i|y_i)P(y_i)}{P(\tilde{\mathbf{x}}_i)}$$

$$P(\tilde{\mathbf{x}}_i) = P(\tilde{\mathbf{x}}_i | y_i = 1) P(y_i = 1) + P(\tilde{\mathbf{x}}_i | y_i = 0) P(y_i = 0)$$



Birches versus Ashes

- The last figure also nicely exemplifies the problem of overlapping classes
- Given brightness level as input, one cannot separate the classes and this problem cannot be solved by a more powerful classifier!
- The only way to solve this issue is to use more features (inputs, sensors); for example one might measure spectral amplitudes at different frequencies, including infrared
- Another problem might be that the brightness detector is unreliable ("noisy labels")

Class-specific Distributions

- To model $P(\tilde{\mathbf{x}}_i|y_i)$ one can chose an application specific distribution
- A popular choice is a Gaussian distribution (normal discriminant analysis)

$$P(\tilde{\mathbf{x}}_i|y_i=l) = \mathcal{N}(\tilde{\mathbf{x}}_i; \vec{\mu}^{(l)}, \boldsymbol{\Sigma})$$

1 - >

with

$$\mathcal{N}\left(\tilde{\mathbf{x}}_{i}; \vec{\mu}^{(l)}, \boldsymbol{\Sigma}\right) = \frac{1}{(2\pi)^{M/2} \sqrt{|\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2} \left(\tilde{\mathbf{x}}_{i} - \vec{\mu}^{(l)}\right)^{T} \boldsymbol{\Sigma}^{-1} \left(\tilde{\mathbf{x}}_{i} - \vec{\mu}^{(l)}\right)\right)$$

 Note, that both Gaussian distributions have different modes (centers) but the same covariance matrices. This has been shown to often work well



Multivariate Normal Dehrbution



Maximum-likelihood Estimators for Modes and Covariances

• One obtains a maximum likelihood estimators for the modes

$$\widehat{\vec{\mu}}^{(l)} = \frac{1}{N_l} \sum_{i: y_i = l} \tilde{\mathbf{x}}_i$$

• One obtains as unbiased estimators for the covariance matrix

$$\widehat{\boldsymbol{\Sigma}} = \frac{1}{N-M} \sum_{l=0}^{1} \sum_{i:y_i=l}^{1} (\widetilde{\mathbf{x}}_i - \widehat{\vec{\mu}}^{(l)}) (\widetilde{\mathbf{x}}_i - \widehat{\vec{\mu}}^{(l)})^T$$

Expanding the Quadratic Terms in the Exponent

• Note that

$$-\frac{1}{2} \left(\tilde{\mathbf{x}}_{i} - \vec{\mu}^{(l)} \right)^{T} \Sigma^{-1} \left(\tilde{\mathbf{x}}_{i} - \vec{\mu}^{(l)} \right)$$
$$= -\frac{1}{2} \tilde{\mathbf{x}}_{i}^{T} \Sigma^{-1} \tilde{\mathbf{x}}_{i} - \frac{1}{2} \vec{\mu}^{(l)}^{T} \Sigma^{-1} \vec{\mu}^{(l)} + \vec{\mu}^{(l)}^{T} \Sigma^{-1} \tilde{\mathbf{x}}_{i}$$

The Difference of the Quadratic

• Now we calculate the difference of the quadratic terms of the two Gaussians

$$-\frac{1}{2} \left(\tilde{\mathbf{x}}_{i} - \vec{\mu}^{(0)} \right)^{T} \Sigma^{-1} \left(\tilde{\mathbf{x}}_{i} - \vec{\mu}^{(0)} \right) + \frac{1}{2} \left(\tilde{\mathbf{x}}_{i} - \vec{\mu}^{(1)} \right)^{T} \Sigma^{-1} \left(\tilde{\mathbf{x}}_{i} - \vec{\mu}^{(1)} \right)$$
$$= -\frac{1}{2} \tilde{\mathbf{x}}_{i}^{T} \Sigma^{-1} \tilde{\mathbf{x}}_{i} - \frac{1}{2} \vec{\mu}^{(0)} \Sigma^{-1} \vec{\mu}^{(0)} + \vec{\mu}^{(0)} \Sigma^{-1} \tilde{\mathbf{x}}_{i}$$
$$+ \frac{1}{2} \tilde{\mathbf{x}}_{i}^{T} \Sigma^{-1} \tilde{\mathbf{x}}_{i} + \frac{1}{2} \vec{\mu}^{(1)} \Sigma^{-1} \vec{\mu}^{(1)} - \vec{\mu}^{(1)} \Sigma^{-1} \tilde{\mathbf{x}}_{i}$$

• since two terms cancel,

$$= \left(\vec{\mu}^{(0)} - \vec{\mu}^{(1)}\right)^T \Sigma^{-1} \tilde{\mathbf{x}}_i - \frac{1}{2} \vec{\mu}^{(0)T} \Sigma^{-1} \vec{\mu}^{(0)} + \frac{1}{2} \vec{\mu}^{(1)T} \Sigma^{-1} \vec{\mu}^{(1)}$$

A Posteriori Distribution

• It follows that

$$P(y_{i} = 1 | \tilde{\mathbf{x}}_{i}) = \frac{P(\tilde{\mathbf{x}}_{i} | y_{i} = 1) P(y_{i} = 1)}{P(\tilde{\mathbf{x}}_{i} | y_{i} = 1) P(y_{i} = 1) + P(\tilde{\mathbf{x}}_{i} | y_{i} = 0) P(y_{i} = 0)}$$

$$= \frac{1}{1 + \frac{P(\tilde{\mathbf{x}}_{i} | y_{i} = 0) P(y_{i} = 0)}{P(\tilde{\mathbf{x}}_{i} | y_{i} = 1) P(y_{i} = 1)}}$$

$$= \frac{1}{1 + \frac{\kappa_{0}}{\kappa_{1}} \exp\left((\vec{\mu}^{(0)} - \vec{\mu}^{(1)})^{T} \Sigma^{-1} \tilde{\mathbf{x}}_{i} - \frac{1}{2} \vec{\mu}^{(0)^{T}} \Sigma^{-1} \vec{\mu}^{(0)} + \frac{1}{2} \vec{\mu}^{(1)^{T}} \Sigma^{-1} \vec{\mu}^{(1)}\right)}$$

$$= \operatorname{sig}\left(w_{0} + \tilde{\mathbf{x}}_{i}^{T} \tilde{\mathbf{w}}\right) = \operatorname{sig}\left(w_{0} + \sum_{j=1}^{M} x_{i,j} w_{j}\right)$$

Weights

• We get ($\tilde{\mathbf{w}}$ is without w_0)

$$\tilde{\mathbf{w}} = \boldsymbol{\Sigma}^{-1} \left(\vec{\mu}^{(1)} - \vec{\mu}^{(0)} \right)$$

- Note that \tilde{w} is independent of κ_1 and κ_0 and is thus independent of the class proportions in the training data! This is important, e.g., for case-control studies
- Recall: sig(arg) = 1/(1 + exp(-arg))

Bias Term

• We get,

$$w_0 = \log \kappa_1 / \kappa_0 + \frac{1}{2} \vec{\mu}^{(0)T} \Sigma^{-1} \vec{\mu}^{(0)} - \frac{1}{2} \vec{\mu}^{(1)T} \Sigma^{-1} \vec{\mu}^{(1)}$$

• w_0 clearly reflects the class proportions



Comments

- This specific generative model leads to linear class boundaries
- But we do not only get class boundaries, we get probabilities
- Although we have used Bayes formula, the analysis was frequentist. A Bayesian analysis with a prior distribution on the parameters is also possible

Comments (cont'd)

- If the two class-specific Gaussians have different covariance matrices $(\Sigma^{(0)}, \Sigma^{(1)})$ the approach is still feasible but one would need to estimate two covariance matrices and the decision boundaries are not linear anymore; still, one can simply apply Bayes rule to obtain posterior probabilities
- The generalization to multiple classes is straightforward: simply estimate a different Gaussian for each class (with shared covariances or not) and apply Bayes rule
- *Generative-Discriminative pair*: (1) Gaussian Analysis (as a generative model) and (2) logistic regression as a discriminant model
- Generalization to basis functions is straight forward: x is replaced by $\vec{\phi}(\mathbf{x})$
- With an explicit $P(\tilde{\mathbf{x}}_i | y_i = l) = \mathcal{N}(\tilde{\mathbf{x}}_i; \vec{\mu}^{(l)}, \Sigma)$, we can apply Bayes formula for a posteriori class estimation
- This is not easy, or even impossible, e.g., for GANs, which are able to generate samples but where the likelihood is not easily evaluated (likelihood free methods)

Special Case: Naive Bayes

• With diagonal covariances matrices, one obtains a *Naive-Bayes* classifier

$$P(\tilde{\mathbf{x}}_i|y_i=l) = \prod_{j=1}^M \mathcal{N}(x_{i,j}; \mu_j^{(l)}, \sigma_j^2)$$

- The naive Bayes classifier has considerable fewer parameters but completely ignores class-specific correlations between features; this is sometimes considered to be naive
- Even more naive (all Gaussian have identical variance):

$$P(\tilde{\mathbf{x}}_i|y_i=l) = \prod_{j=1}^M \mathcal{N}(x_{i,j}; \mu_j^{(l)}, \sigma^2)$$

Logistic Regression from Naive Bayes

• We have parameters, for the latter case,

$$w_{j} = \frac{1}{\sigma^{2}} \left(\mu_{j}^{(1)} - \mu_{j}^{(0)} \right)$$
$$w_{0} = \log \kappa_{1} / \kappa_{0} + \frac{1}{2\sigma^{2}} \sum_{j} \left(\mu_{j}^{(0)} \right)^{2} - \left(\mu_{j}^{(1)} \right)^{2}$$

- Note that w_j is completely independent of other inputs; adding or removing other inputs does not change w_j;
- In contrast w_0 depends on all dimensions
- The smaller σ^2 , the sharper the transition

Special Case: Bernoulli Naive Bayes

- Naive Bayes classifiers are popular in text analysis with often more than 10000 features (key words). For example, the classes might be SPAM and no-SPAM and the features are keywords in the texts
- Instead of a Gaussian distribution, a Bernoulli distribution is employed
- $P(word_j = 1 | SPAM) = \gamma_{j,s}$ is the probability of observing word $word_j$ in the document for SPAM documents (Bernoulli distribution)

Special Case: Bernoulli Naive Bayes

- We also consider the other cases
- $P(word_j = 0|SPAM) = 1 \gamma_{j,s}$ is the probability of not observing word $word_j$ in the document for SPAM documents
- $P(word_j = 1 | no-SPAM) = \gamma_{j,n}$ is the probability of observing word $word_j$ in the document for non-SPAM documents
- $P(word_j = 0|no-SPAM) = 1 \gamma_{j,n}$ is the probability of not observing word $word_j$ in the document for non-SPAM documents
- Note that there are two parameters per dimension: $\gamma_{j,s}$ and $\gamma_{j,n}$

Special Case: Bernoulli Naive Bayes (cont'd)

• Then

P(SPAM|doc) =

$$\frac{\kappa_s \prod_j \gamma_{j,s}^{\mathsf{word}_j} \ (1 - \gamma_{j,s})^{1 - \mathsf{word}_j}}{\kappa_s \prod_j \gamma_{j,s}^{\mathsf{word}_j} \ (1 - \gamma_{j,s})^{1 - \mathsf{word}_j} + \kappa_n \prod_j \gamma_{j,n}^{\mathsf{word}_j} \ (1 - \gamma_{j,n})^{1 - \mathsf{word}_j}}$$

• Simple ML estimates are $\gamma_{j,s} = N_{j,s}/N_s$ and $\gamma_{j,n} = N_{j,n}/N_n$

(N_s is the number of SPAM documents in the training set, $N_{j,s}$ is the number of SPAM documents in the training set where *word*_j is present)

(N_n is the number of no-SPAM documents in the training set, $N_{j,n}$ is the number of no-SPAM documents in the training set where *word*_j is present)

Special Case: Bernoulli Naive Bayes (cont'd)

• Note, that we can also write

$$P(\mathsf{SPAM}|doc) = sig(w_0 + \sum_j w_j word_j)$$

with

$$w_j = [\log \gamma_{j,s} - \log \gamma_{j,n}] - [\log(1 - \gamma_{j,s}) - \log(1 - \gamma_{j,n})]$$
$$w_0 = \log \kappa_s / \kappa_n + \sum_j \log(1 - \gamma_{j,s}) - \log(1 - \gamma_{j,n})$$

• Generative-Discriminative pair: (1) Bernoulli naive Bayes classifier and (2) logistic regression

Probabilistic Mixture Models

- With K classes, we obtain the joint distribution $P(Y = k)P(\tilde{\mathbf{x}}|Y = k)$, where $k = 1, \dots, K$
- Consider the generative models just discussed but assume that the class labels are unknown during training
- Then we achieve probabilistic mixture models: Mixture of Gaussians, Mixture of Multinomials, Mixture of Bernoullis, ...
- One discovers hidden classes or clusters: For example, the centres of the Gaussians after training can be interpreted as cluster representatives
- At the same time, we obtain a model for the data distribution $P(\tilde{\mathbf{x}}) = \sum_{k=1}^{K} P(\tilde{\mathbf{x}}|Y=k)P(Y=k)$
- Training can be done using the Expectation Maximization (EM) algorithm

II. Logistic Regression

- In I. (Generative models for classification) we first defined a generative model for P(x, y); from this model we then derived P(y|x) = P(y)P(x|y) which models x given y (generative modelling)
- Here, we model the reverse $P(y|\mathbf{x})$ (standard supervised learning)
- With logistic regression as the discriminant version, we model discriminatively

$$\hat{y}_i = P(y = 1 | \mathbf{x}_i) = \operatorname{sig}\left(\mathbf{x}_i^T \mathbf{w}\right)$$

(now we include the bias $\mathbf{x}_i^T = (x_{i,0} = 1, x_{i,1}, \dots, x_{i,M-1})^T$). Sig() as defined before (logistic function).

• One now optimizes the likelihood of the conditional model

$$L(\mathbf{w}) = \prod_{i=1}^{N} \operatorname{sig} \left(\mathbf{x}_{i}^{T} \mathbf{w} \right)^{y_{i}} \left(1 - \operatorname{sig} \left(\mathbf{x}_{i}^{T} \mathbf{w} \right) \right)^{1-y_{i}}$$

Log-Likelihood Function is the Negative Cross Entropy

• Log-likelihood function

$$l = \sum_{i=1}^{N} y_i \log \left(\operatorname{sig} \left(\mathbf{x}_i^T \mathbf{w} \right) \right) + (1 - y_i) \log \left(1 - \operatorname{sig} \left(\mathbf{x}_i^T \mathbf{w} \right) \right)$$

• Cross-entropy cost function (negative log-likelihood)

$$l = -\left[\sum_{i=1}^{N} y_i \log\left(\operatorname{sig}\left(\mathbf{x}_i^T \mathbf{w}\right)\right) + (1 - y_i) \log\left(1 - \operatorname{sig}\left(\mathbf{x}_i^T \mathbf{w}\right)\right)\right]$$

Log-Likelihood

• Log-likelihood function

$$l = \sum_{i=1}^{N} y_i \log\left(\operatorname{sig}\left(\mathbf{x}_i^T \mathbf{w}\right)\right) + (1 - y_i) \log\left(1 - \operatorname{sig}\left(\mathbf{x}_i^T \mathbf{w}\right)\right)$$
$$l = \sum_{i=1}^{N} y_i \log\left(\frac{1}{1 + \exp(-\mathbf{x}_i^T \mathbf{w})}\right) + (1 - y_i) \log\left(\frac{1}{1 + \exp(\mathbf{x}_i^T \mathbf{w})}\right)$$

$$= -\sum_{i=1}^{N} y_i \log(1 + \exp(-\mathbf{x}_i^T \mathbf{w})) + (1 - y_i) \log(1 + \exp(\mathbf{x}_i^T \mathbf{w}))$$

Adaption

• The derivatives of the log-likelihood with respect to the parameters

$$\frac{\partial l}{\partial \mathbf{w}} = \sum_{i=1}^{N} y_i \frac{\mathbf{x}_i \exp(-\mathbf{x}_i^T \mathbf{w})}{1 + \exp(-\mathbf{x}_i^T \mathbf{w})} - (1 - y_i) \frac{\mathbf{x}_i \exp(\mathbf{x}_i^T \mathbf{w})}{1 + \exp(\mathbf{x}_i^T \mathbf{w})}$$
$$= \sum_{i=1}^{N} y_i \mathbf{x}_i (1 - \operatorname{sig}(\mathbf{x}_i^T \mathbf{w})) - (1 - y_i) \mathbf{x}_i \operatorname{sig}(\mathbf{x}_i^T \mathbf{w})$$

$$= \sum_{i=1}^{N} (y_i - \operatorname{sig}(\mathbf{x}_i^T \mathbf{w})) \mathbf{x}_i = \sum_{i=1}^{N} (y_i - \hat{y}_i) \mathbf{x}_i$$

• A gradient-based optimization of the parameters to maximize the log-likelihood

$$\mathbf{w} \longleftarrow \mathbf{w} + \eta \frac{\partial l}{\partial \mathbf{w}}$$

• Typically one uses a Newton-Raphson optimization procedure



Logistic Regression as a Generalized Linear Models (GLM)

- Consider a Bernoulli distribution with $P(y = 1) = \theta$ and $P(y = 0) = 1 \theta$, with $0 \le \theta \le 1$
- In the theory of the exponential family of distributions, one sets

 $\theta = \operatorname{sig}(\eta)$

Now we get valid probabilities for any $\eta \in \mathbb{R}!$

• η is called the natural parameter and Sig(·) the inverse parameter mapping for the Bernoulli distribution

Logistic Regression as a Generalized Linear Models (GLM) (cont'd)

• This is convenient if we make η a linear function of the inputs and one obtains a Generalized Linear Model (GLM)

$$P(y_i = 1 | \mathbf{x}_i, \mathbf{w}) = sig(\mathbf{x}_i^T \mathbf{w})$$

• Thus logistic regression is the GLM for the Bernoulli likelihood model

Application to Neural Networks and other Systems

- Logistic regression essentially defines a new cost function
- It can be applied as well to neural networks, as we have done before,

$$P(y_i = 1 | \mathbf{x}_i, \mathbf{w}) = sig(NN(\mathbf{x}_i))$$

or systems of basis functions or kernel systems

Multiple Classes and Softmax

- Consider a multinomial distribution with $P(y = c) = \theta_c$, with $\theta_c \ge 0$ and $\sum_{c=1}^{C} \theta_c = 1$. c is the class index and C is the number of classes
- We reparameterize (exponential family of distributions)

$$\theta_c = \frac{\exp(\eta_c)}{\sum_{c'=1}^{C} \exp(\eta_{c'})}$$

• The η_c are unconstrained; softmax notation: $\theta_c = \operatorname{softmax}_c(\vec{\eta_c})$

Multiple Classes and Softmax: GLM

$$ullet$$
 In GLM, we set $\eta_c = \mathbf{x}^T \mathbf{w}_c$ and

$$\hat{y}_c = P(y = c | \mathbf{x}) = \frac{\exp(\mathbf{x}^T \mathbf{w}_c)}{\sum_{c'=1}^{C} \exp(\mathbf{x}^T \mathbf{w}_{c'})}$$

Multiple Classes and Softmax (cont'd)

• The negative log-likelihood (softmax cross entropy) becomes

$$-l = -\sum_{i=1}^{N} \left(\sum_{c=1}^{C} y_{i,c} \mathbf{x}_{i}^{T} \mathbf{w}_{c} - \log \sum_{c=1}^{C} \exp(\mathbf{x}_{i}^{T} \mathbf{w}_{c}) \right)$$

Multiple Classes and Softmax (cont'd)

• The gradient becomes

$$-\frac{\partial l}{\partial w_{j,c}} = -\sum_{i} \left(y_{i,c} x_{i,j} - \frac{x_{i,j} \exp(\mathbf{x}_{i}^{T} \mathbf{w}_{c})}{\sum_{c=1}^{C} \exp(\mathbf{x}_{i}^{T} \mathbf{w}_{c})} \right)$$

and SGD becomes

$$w_{j,c} \leftarrow w_{j,c} + \eta x_{i,j} (y_{i,c} - \hat{y}_{i,c})$$

III. Classification via Regression

• Linear Regression:

$$f(\mathbf{x}_i, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j x_{i,j}$$
$$= \mathbf{x}_i^T \mathbf{w}$$

- We define as target $y_i = 1$ if the pattern \mathbf{x}_i belongs to class 1 and $y_i = 0$ (or $y_i = -1$) if pattern \mathbf{x}_i belongs to class 0
- We calculate weights $\mathbf{w}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ as LS solution, exactly as in linear regression
- For a new pattern x we calculate $f(\mathbf{x}) = \mathbf{x}^T \mathbf{w}_{LS}$ and assign the pattern to class 1 if $f(\mathbf{x}) > 1/2$ (or $f(\mathbf{x}) > 0$); otherwise we assign the pattern to class 0

Bias

- Asymptotically, a LS-solution converges to the posterior class probabilities, although a linear functions is typically not able to represent $P(c = 1|\mathbf{x})$. The resulting class boundary can still be sensible
- One can expect good class boundaries in high dimensions and/or in combination with basis functions, kernels and neural networks; in high dimensions sometimes consistency can be achieved. In essence it is necessary that the linear model can model the expected probability $P(c = 1|\mathbf{x})$

Classification via Regression with Linear Functions



Classification via Regression with Radial Basis Functions



Causal Effect

- Assume that all relevant inputs are considered in the model (no other confounders) and that we use "Classification via Regression"
- The causal effect is independent of the individual, and can be estimated as

$$P(Y = 1 | x_{i,1} = 1, x_{i,2}, \dots, x_{i,M}) - P(Y = 1 | x_{i,1} = 0, x_{i,2}, \dots, x_{i,M}) = w_1$$

- $x_1 = 1$ means that the individual has received the treatment, and $x_1 = 0$ means that the individual has not received the treatment,
- Y = 1 means that the patient is healthy after the treatment

Performance

- Although the approach might seem simplistic, the performance can be excellent (in particular in high dimensions and/or in combination with basis functions, kernels and neural networks). The calculation of the optimal parameters can be very fast!
- Regression is commonly used in treatment effect prediction in medicine if the influence of the treatment is small, on average

Logistic Regression in Medical Statistics

Logistic Regression in Medical Statistics

- Logistic regression has become one of the the most important tools in medical statistics to analyse the outcome of treatments, e.g., a new medication, and to evaluate the effect of preconditions (gender, age, smoking, environmental effects)
- An important task is to distinguish between correlation and causation

Epidemiology

- In epidemiology, the output y = 1 means that the patient has the disease
- $x_1 = 1$ might represent the fact that the patient was exposed (e.g., by a genetic variant, smoking, or an environmental factor) and $x_1 = 0$ might mean that the patient was not exposed; the other inputs are often typical confounders (age, sex, ...)

$$P(y_i = 1 | \mathbf{x}_i, \mathbf{w}) = \operatorname{sig} \left(\sum_{j=0}^{M} w_j x_j \right)$$

- Thus, w_1 is the quantity of interest! If w_1 is significantly larger than zero, then the exposure was harmful!
- For model fitting we need data from individuals, which were randomly chosen out off the population; for rare diseases, his can be a problem (see later discussion on the logs-odds ratio)

Treatment Evaluation

- All individuals in the population have the disease
- In treatment evaluation, $x_1 = 1$ means that the patient received the treatment, and $x_1 = 0$ means that the patient did not receive the treatment
- The output represents the outcome after treatment; e.g., y = 1 can mean that the patient is cured by the treatment

$$P(y_i = 1 | \mathbf{x}_i, \mathbf{w}) = \operatorname{sig} \left(\sum_{j=0}^{M} w_j x_j \right)$$

• Of course, of great interest is if w_1 is significantly nonzero

Causal Effect Depends on the Individual

• In the model, the causal effect depends on the individual,

$$P(y = 1 | x_{i,1} = 1, x_{i,2}, \dots, x_{i,M}) - P(y = 1 | x_{i,1} = 0, x_{i,2}, \dots, x_{i,M})$$
$$= sig(x_{i,1} = 1, x_{i,2}, \dots, x_{i,M}) - sig(x_{i,1} = 0, x_{i,2}, \dots, x_{i,M})$$

- We can calculate the average causal effect, e.g., on subgroups (stratification))
- Maybe we can also find an interpretation of w_1 , which we analyse next

Odds

• The **odds** for a patient with properties \mathbf{x}_i is defined as

$$Odds(\mathbf{x}_i) = \frac{P(y_i = 1 | \mathbf{x}_i)}{P(y_i = 0 | \mathbf{x}_i)}$$

Log-Odds for Logistic Regression

- In medical statistics, one is interested in the interpretation of the terms in logistic regression
- For logistic regression, the log odds (= logit) is

$$LogOdds = \log \frac{P(y_i = 1 | \mathbf{x}_i)}{P(y_i = 0 | \mathbf{x}_i)} = \log \frac{1}{1 + \exp(-\mathbf{x}_i^T \mathbf{w})} \frac{1 + \exp(-\mathbf{x}_i^T \mathbf{w})}{\exp(-\mathbf{x}_i^T \mathbf{w})}$$
$$= \log \frac{1}{\exp(-\mathbf{x}_i^T \mathbf{w})} = \mathbf{x}_i^T \mathbf{w}$$

- Thus the log odds of the outcome is $h = \mathbf{x}_i^T \mathbf{w}$, which is the net input, also called the logit (with $y = \operatorname{sig}(x)$, $x = \operatorname{sig}^{-1}(y) = \log(y) - \log(1 - y)$)
- Thus logistic regression is linear in predicting the log odds

(Log) Odds Ratio

• The odds ratio is defined as

$$OR = \frac{Odds(x_{i,1} = 1, x_{i,2}, \dots, x_{i,M})}{Odds(x_{i,1} = 0, x_{i,2}, \dots, x_{i,M})}$$

• The log odds ratio evaluates the effect of the treatment

 $\log(OR) = \log Odds(x_{i,1} = 1, x_{i,2}, \dots, x_{i,M}) - \log Odds(x_{i,1} = 0, x_{i,2}, \dots, x_{i,M})$

The Log Odds Ratio for Logistic Regression

• In logistic regression, the log odds ratio is identical to w_1 , since

$$(w_0 + w_1 + \sum_{j=2}^M x_{i,j}) - (w_0 + 0 + \sum_{j=2}^M x_{i,j}) = w_1$$

- If w_1 is significantly nonzero, then the exposure/treatment has an effect
- Thus, in logistic regression, the causal effect might be different for each individual, the log odds ratio is the same; $w_1 \ge 0$ is the increase in the log odds of a patient that obtains the treatment, independent of the confounders $x_{i,2}, \ldots, x_{i,M}$)
- If possible, confounders (x_{i,2},..., x_{i,M}) should be inputs to the logistic regression (e.g., age, gender) trained on the population; this is called: "controlling for the confounders" or "controlling confounding effects"
- Alternatively, one forms subgroups of patients with the same values of the confounders (same age group, same gender) and calculates the log odds ratio separately for each subgroup (stratification) (so the logistic regression input is only $x_{i,1}$)

Case Control Studies and Imbalanced Classes

- Consider a rare disease that only affects one in a million; then if I would collect data from 1 million random individuals I might only have one individual with the disease
- Applying Bayes formula, the odds ratio can also be written as

$$OR = \frac{P(x_{i,1} = 1 | y_i = 1, x_{i,2}, \dots, x_{i,M})}{P(x_{i,1} = 0 | y_i = 1, x_{i,2}, \dots, x_{i,M})} \frac{P(x_{i,1} = 0 | y_i = 0, x_{i,2}, \dots, x_{i,M})}{P(x_{i,1} = 1 | y_i = 0, x_{i,2}, \dots, x_{i,M})}$$

- To obtain these conditional probabilities we can select patients according to their disease status and any other attributes (but then not according to the fact if they got exposed or not)
- It is the quite realistic to collect let's say 1000 individuals who have the disease (e.g., breast cancer) and 1000 who do not have the disease, and then predict if they were exposed or not!
- In a model that predicts treatment from outcome where x₁ is the outcome, w₁ again is the log odds ratio

Odds Ration and Relative Risk

• The relative risk (= risk ratio) (*RR*) is considered more intuitive (see English Wikipedia page on "odds ratio"), but is does not possess this insensitivity to sampling

$$RR = \frac{P(y_i = 1 | x_{i,1} = 1, x_{i,2}, \dots, x_{i,M})}{P(y_i = 1 | x_{i,1} = 0, x_{i,2}, \dots, x_{i,M})}$$

• For rare diseases (outcomes) *RR* and *OR* are very similar

Causality

- Confounders are variables that influence both the output y and x_1
- For example R. A. Fisher argued that there might be a genetic variant which makes you want to smoke and which gives you lung cancer; thus you would get lung cancer independently if you smoked
- This turned out to be (mostly) untrue
- By far the most apparent paradoxes result from unmodelled confounders (Simpson's paradox)

Linear Regression versus Logistic Regression

• If I apply linear regression, the *causal effect* (*CE*) is independent of the individual, and can be estimated as

$$w_1^{\text{linear regression}} = CE$$

$$= P(Y = 1 | x_{i,1} = 1, x_{i,2}, \dots, x_{i,M}) - P(Y = 1 | x_{i,1} = 0, x_{i,2}, \dots, x_{i,M})$$

• If I apply logistic regression,

$$w_1^{\text{logistic regression}} = \log OR$$

• Relative risk:

 $\log RR =$

 $\log P(Y = 1 | x_{i,1} = 1, x_{i,2}, \dots, x_{i,M}) - \log P(Y = 1 | x_{i,1} = 0, x_{i,2}, \dots, x_{i,M})$





Rare Diseases: Random Sample

Random individuals from the population

Disease	2	3
No Disease	100	100
	not exposed	exposed

- $OR_{y|x} = (3/103)/(100/103) \times (100/102)/(2/102) = 1.5$
- CE = 3/(3+100) 2/(2+100) = 0.0095
- RR = 3/(3+100)/(2/(2+100)) = 0.149
- $OR_{x|y} = (3/5)/(2/5) \times (100/200)/(100/200) = 1.5$

Rare Diseases: Balanced

(Disease / No disease) balance

Disease	80	120
No Disease	100	100
	not exposed	exposed

- $OR_{y|x} = (120/220)/(100/220) \times (100/180)/(80/180) = 1.5$
- $CE = \frac{120}{(120 + 100)} \frac{80}{(80 + 100)} = 0.101$
- $RR = \frac{120}{(120 + 100)} / (\frac{80}{(80 + 100)}) = 1.227$
- $OR_{x|y} = (120/200)/(80/200) \times (100/200)/(100/200) = 1.5$
- $OR_{y|x}$ and $OR_{x|y}$ remain at 1.5, but can be estimated with greater accuracy

The Rubin causal model (RCM),

- Select a patient i who has received the treatment $x_{1,i} = 1$
- Find a twin who did not receive the treatment: Select another patient c(i) who is as identical as possible to the first patient, according to the values of the confounders $x_{2,i}, \ldots, x_{M,i}$, but who is in the control group, i.e., who did not receive the treatment $(x_{1,c(i)} = 0)$
- Repeat this for all N patients who received the treatment
- This is a counterfactual analysis; the causal effect is estimated as

$$CE = \frac{1}{N} \sum_{i=1}^{N} (y_i - y_{c(i)}) = \frac{1}{N} \sum_{i=1}^{N} y_i - \frac{1}{N} \sum_{i=1}^{N} y_{c(i)}$$

Personalized Medicine

- A linear model assumes that the effect of an input on the output is independent of the other inputs
- A log-linear model assumes that the effect of an input on the log-odds of the output is independent of the other inputs
- The idea behind personalized medicine is that a given medication only works for a subclass of the population
- Thus one either tries to identify as good as possible the group (strata) for which the medication works or one includes nonlinear interactions between the inputs
- If many factors might contribute to the effectiveness of a drug, one might try multivariate nonlinear models, e.g., neural networks