

Supervised Classification



Setup: Categorical Label Prediction

Classification tasks aim at predicting a discrete output $y \in \mathcal{Y} = \{C_1, \dots, C_g\}$ with $2 \leq g < \infty$ given Data D .

We encode \mathcal{Y} as: binary case $g=2$: $\mathcal{Y} = \{0, 1\}$ or $\mathcal{Y} = \{-1, 1\}$.

usually ordered classes: multiclass case $g>2$: $\mathcal{Y} = \{1, \dots, g\}$.

Our model $f: \mathcal{X} \rightarrow \mathbb{R}^g$ outputs scores/probabilities and not classes.

- Continuous function is easier to optimize than discrete-valued functions.
- Scores/Probabilities contain more information and can be transformed into classes. But this transformation is non-injective.

Approaches to construct classifiers

Goal: Model joint data-generating process to obtain class prob directly or being in class k given feature vector x .

Generative Approach: Model posterior probability $\pi_k(x)$ using Bayes

Idea: Assume that the data generating process for x depends on y .

Then we want to answer "Which y tends to have x like that?"

$$\pi_k(x) = P(y=k|x) = \frac{P(x|y=k) \cdot \pi_k}{\sum_{j=1}^g P(x|y=j) \cdot \pi_j}$$

with π_j being a prior estimated from data, e.g. using e.g. relative frequencies.

Model posterior distributions, then choose best decision boundary based on that

Logistic (sigmoid) function: $P(y=1|x, \theta) = \pi(x|\theta)$

for binary classifier: $S(\cdot)$ estimates prob. of class "1"

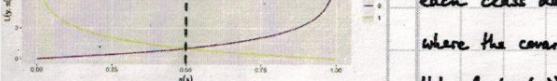
$$H = \{x \in \mathcal{X} \mid \pi(x) \geq \frac{1}{2}\}$$

$$s(\theta^T x) = \frac{\exp(\theta^T x)}{1 + \exp(\theta^T x)}$$

Logistic / Bernoulli Loss function:

$$-(y, \pi(x)) = -y \log(\pi(x)) + (1-y) \log(1 - \pi(x))$$

Optimization is done numerically (GD).



Analyses heavier the more confident we are in a wrong prediction! (Loss $\rightarrow \infty$)

Usually we don't directly predict hard labels $h(\cdot)$. Scores / Probabilities are more informative and it's much easier to work with continuous values in optimization

Scoring classifier

Construct g discriminant / scoring functions $f_1, \dots, f_g: \mathcal{X} \rightarrow \mathbb{R}$.

We choose the class with the maximum score, i.e. $h(x) = \arg \max_{k=1, \dots, g} f_k(x)$.

Special case $g=2$: $f_1(x)$ is "preference" score for positive class (always 1 iff).

One discriminant / scoring function is sufficient: $f(x) = f_1(x) - f_2(x)$ (Note: $\mathcal{Y} = \{-1, 1\}$).

$h(x) = \text{sgn}(f(x))$ and $|f(x)|$ is called confidence.

Linear Classifier: Affine linear in features x . If we allow transf. of feature space (e.g. polynomial) we can handle non-linear classif. in original space (can learn non-linear boundary).

If the discriminant / scoring functions f_1, \dots, f_g can be specified as linear functions, i.e. $f_k(x) = w_k^T x + b_k$ with g being a rank-preserving, monotone transformation, then we call f_1, \dots, f_g linear classifiers.

The decision boundaries are hyperplanes.

$$f_k(x) = w_k^T x + b_k \Rightarrow w_k^T x + b_k = 0$$

Binary case: Convert scores or prob to class outputs by thresholding $h(x) = \mathbb{1}[f(x) \geq c]$; standard $c=0$ for scores, $c=0.5$ for probs.

Discriminant Approach: Goal: Directly learn Decision boundary with ERM

Idea: Use empirical risk minimization with a suitable Loss-function.

$$\hat{f} = \arg \min_{f \in \mathcal{F}} R_{\text{emp}}(f) = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n L(f(x_i), y_i)$$

We only have limited data n , if p gets too large for that then we may prefer simpler LDA approach.

Reminder: LDA might be preferable over QDA in higher dimensions of features $p \gg n$.

Only sensible for numerical features.

Multivariate gaussian distribution in each class.

$$p(x|y=k) = \frac{1}{(2\pi)^{p/2} |\Sigma_k|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)\right)$$

Linear discriminant analysis assumes that each class density is modeled as $x|y=k \sim \mathcal{N}(\mu_k, \Sigma)$.

where the covariance Σ is equal for all classes.

Note: $\pi_k(x) \propto p(x|y=k) \cdot \pi_k \propto \exp(\theta_k^T x)$.

Linear decision boundary $\Rightarrow \log(\pi_k(x)) = \theta_k^T x + \log(\pi_k)$ LDA gives a linear classifier.

Quadratic discriminant analysis assumes that each class density is modeled as $x|y=k \sim \mathcal{N}(\mu_k, \Sigma_k)$.

where the covariances Σ_k do not need to be equal.

Note: $\log(\pi_k(x))$ defines a quadratic function in x , i.e. not linear.

\Rightarrow quadratic decision boundary. But if quadratic terms cancel out we can also get linear boundary!

Probabilistic classifier

Construct g probability functions $\pi_1, \dots, \pi_g: \mathcal{X} \rightarrow [0, 1]$ with $\sum_{k=1}^g \pi_k = 1$.

We choose the class with the maximum score, i.e. $h(x) = \arg \max_{k=1, \dots, g} \pi_k(x)$.

Special case $g=2$: One discriminant / scoring function is sufficient: $\pi(x) = 1 - \pi_2(x)$ (Note: $\mathcal{Y} = \{0, 1\}$).

$h(x) = \mathbb{1}[\pi(x) \geq c]$ with c being a specified threshold.

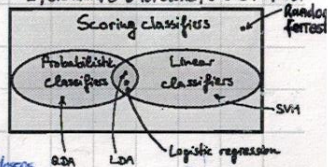
Decision Region

A decision region for class k is defined as $X_k = \{x \in \mathcal{X} \mid h(x) = k\}$.

The decision boundaries are defined as hypersurfaces in \mathcal{X} with tied maximal score:

$$X_k = \{x \in \mathcal{X} \mid \exists i, j \in \{1, \dots, g\}: f_i(x) = f_j(x) \wedge \forall k, i, j: f_i(x) \geq f_k(x) \wedge f_j(x) \geq f_k(x)\}$$

Special case $g=2$ with threshold c : the decision boundary is $\{x \in \mathcal{X} \mid f(x) = c\}$.



Connections and differences of LDA, QDA and NB

Assuming numerical features, cause LDA, QDA require that!

Fill model conditional distr. of features $p(x|y=k)$ as multivariate Gaussian.

QDA: $x|y=k \sim \mathcal{N}(\mu_k, \Sigma_k)$ LDA: $x|y=k \sim \mathcal{N}(\mu_k, \Sigma)$ $\Sigma_k = \Sigma$ same cov. for all classes.

LDA and NB are special cases! (GNB: $x|y=k \sim \mathcal{N}(\mu_k, D_k)$ $\Sigma_k = D_k = \text{diag}(\sigma_k^2)$).

Why? NB requires the least amount of parameter estimators for smallish k and large p .

Often NB is the best model of the three although the cond. independence assumption.

Features can be numerical / categorical or mixture!

Naive Bayes (Generative) Idea: $\pi_k(x) \approx P(y=k|x) = \frac{P(x|y=k) \cdot \pi_k}{P(x)}$.

For numerical features: We further assume that $x_j|y=k \sim \mathcal{N}(\mu_{jk}, \sigma_{jk}^2)$.

This results in $x|y=k \sim \mathcal{N}(\mu_k, D_k)$ with D_k being a diagonal covariance matrix.

$$p(x|y=k) = \frac{1}{(2\pi)^{p/2} \prod_{j=1}^p \sigma_{jk}} \exp\left(-\frac{1}{2} \sum_{j=1}^p \frac{(x_j - \mu_{jk})^2}{\sigma_{jk}^2}\right)$$

GNB: different cov. \Rightarrow ellipses, features uncorrelated.

For categorical features: We further assume that $x_j|y=k$ is multinomial distributed, i.e. $p(x_j|y=k) \propto \frac{n!}{x_j! (n-x_j)!} \pi_{jk}^{x_j} (1-\pi_{jk})^{n-x_j}$.

with $n_j = \sum_{k=1}^g x_{jk}$ being the (absolute) / relative frequency of category j in feature x restricted to class k .

Note: The multinomial Naive Bayes classifier is a linear classifier.

joint cond. distr. (prob. for discrete x) = product of factors across features.

cond. independence assumption: $p(x|y=k) = \prod_{j=1}^p p(x_j|y=k)$.

Priors easily estimated as rel. frequencies from training data.

Gaussian Naive Bayes (GNB): $p(x|y=k) = \prod_{j=1}^p p(x_j|y=k)$.

We could choose either distrib. that only need to estimate univariate distr.

GNB is a special case of QDA. $p(x|y=k) = \prod_{j=1}^p \frac{1}{\sigma_{jk}} \exp\left(-\frac{(x_j - \mu_{jk})^2}{2\sigma_{jk}^2}\right)$.

Product of univariate gaussians.

count of x_j in class k .

total count in class k .

cov. $p(x|y=k) = \prod_{j=1}^p p(x_j|y=k)$ follow Gauss?

if x_j is categorical c = number of categories.

piece of category in class k .

global histogram.

with cond. small $\sigma_{jk} > 0$.

Exercise Classification: Generative Approaches

1.] Proof that LDA is a linear classifier

(feature) model class distribution as $x|y=k \sim \mathcal{N}(M_k, \Sigma)$

NB assumes $P(x|y=k) = \prod_{j=1}^p P(x_j|y=k)$ feature independence

Starting point is modelling posterior probability:

of being in class k given feature vector x

$$\pi_k(x) = P(y=k|x) = \frac{\pi_k \cdot p(x|y=k)}{p(x)}$$

proportional to $\pi_k \cdot p(x|y=k)$

proportional; insert $x|y=k \sim \mathcal{N}(M_k, \Sigma)$ without normalise constant

$$\propto \pi_k \cdot p(x|y=k) \propto \pi_k \exp\left(-\frac{1}{2}(x-M_k)^T \Sigma^{-1}(x-M_k)\right) = \pi_k \exp\left(-\frac{1}{2}x^T \Sigma^{-1}x + \frac{1}{2}M_k^T \Sigma^{-1}x - \frac{1}{2}M_k^T \Sigma^{-1}M_k\right)$$

quadratic form in x resulting in scalar (dep. of x) some const. scalar (indep. of x)

$$= \exp(\log \pi_k - \frac{1}{2}M_k^T \Sigma^{-1}M_k + x^T \Sigma^{-1}M_k) \exp(-\frac{1}{2}x^T \Sigma^{-1}x)$$

but independent of class k since Σ is by assumption identical for all classes!

$$= \exp(b_k + x^T w_k) \exp(-\frac{1}{2}x^T \Sigma^{-1}x) \propto \exp(w_k + x^T w_k)$$

Now with $\log()$ transformation we get $\log(\pi_k w_k) \propto w_{0k} + x^T w_k$

$$\text{by defining } w_{0k} := \log \pi_k - \frac{1}{2}M_k^T \Sigma^{-1}M_k \text{ and } w_k := \Sigma^{-1}M_k$$

would still work with scalar $\exp(-\frac{1}{2}x^T \Sigma^{-1}x) > 0$.

transformed score f_k is linear here w_{0k}, w_k are params θ_{0k}, θ_k ignoring some multipl. constants

Contrast this with QDA Model class (feature) distribution as $x|y=k \sim \mathcal{N}(M_k, \Sigma_k)$

ignored $|\Sigma|^{-1/2}$ for LDA as a constant, but here it depends on class $k \Rightarrow$ therefore no further simplifications possible

$$\pi_k(x) \propto \pi_k |\Sigma_k|^{-1/2} \exp\left(-\frac{1}{2}x^T \Sigma_k^{-1}x - \frac{1}{2}M_k^T \Sigma_k^{-1}M_k + x^T \Sigma_k^{-1}M_k\right)$$

(Two-class Problem)

quadratic decision boundary: Assume indifference betw. classes a and b

$$\log(\pi_k) = \log \pi_k - \frac{1}{2} \log |\Sigma_k| - \frac{1}{2}M_k^T \Sigma_k^{-1}M_k + x^T \Sigma_k^{-1}M_k - \frac{1}{2}x^T \Sigma_k^{-1}x$$

$$\log(\pi_a(x)) = \log(\pi_b(x)) \Leftrightarrow -\frac{1}{2}x^T (\Sigma_a^{-1} - \Sigma_b^{-1})x + x^T (\Sigma_a^{-1}M_a - \Sigma_b^{-1}M_b) + b = 0$$

Same for naive Bayes with numerical features Just replace Σ_k with $\text{diag}(\sigma_k^2)$

const. scalar "linear" in x (depending on k)

Quadratic form in x

If $\Sigma_a^{-1} = \Sigma_b^{-1}$ this simplifies to $x^T (\Sigma_a^{-1}(M_a - M_b)) + b = 0$ for same Σ across classes quadratic form disappears & we get linear boundary!

Parameter Estimators in LDA and QDA

LDA: $\hat{\pi}_{ik} = \frac{n_{ik}}{n}$ where n_{ik} is number of class k -observations

QDA: $\hat{\pi}_{ik} = \frac{n_{ik}}{n}$ prior probabilities $\hat{\pi}_{ik}$ class frequencies

$$\hat{M}_k = \frac{1}{n_k} \sum_{i: y_i=k} x^{(i)}$$

$$\hat{M}_k = \frac{1}{n_k} \sum_{i: y_i=k} x^{(i)}$$

$$\hat{\Sigma}_k = \frac{1}{n_k - 1} \sum_{i: y_i=k} (x^{(i)} - \hat{M}_k)(x^{(i)} - \hat{M}_k)^T$$

$$\hat{\Sigma}_k = \frac{1}{n_k - 1} \sum_{i: y_i=k} (x^{(i)} - \hat{M}_k)(x^{(i)} - \hat{M}_k)^T$$

average everything across all classes to get same cov. matrix for all classes

$$\hat{\Sigma} = \frac{1}{n - g} \sum_{k=1}^g \sum_{i: y_i=k} (x^{(i)} - \hat{M}_k)(x^{(i)} - \hat{M}_k)^T$$

Overall: $(k-1) + k \cdot p + \frac{p(p+1)}{2}$ parameters to estimate

Gaussian Naive Bayes

different cov. in each class

$$(GNB: (k-1) + k \cdot p + k \cdot p = (k-1) + 2kp$$

cov. is Diagonal Matrix

Overall: $(k-1) + k \cdot p + \frac{p(p+1)}{2}$ parameters to estimate

p : features

k : classes, may also use g for that

one for each class $\hat{\pi}_{ik}$ \hat{M}_{ks} $\hat{\Sigma}$ p -times on top as \hat{M}_{ks} contains p -averages (one for each feature) $\hat{\Sigma}$ cov. matrix symmetric inner feature variances on diagonal

Note] For QDA we have to estimate way more parameters because of estimation of separate cov. matrix per class

\Rightarrow Weird side effect: LDA might perform better than QDA in higher dimensions even if data-generating process matches QDA assumptions (fully or better)!!

class error / cross-validation

in terms of feature space, e.g. 100+ features ($p \approx 100$)

The Reason is that we have a limited amount of data (n) and therefore we can't estimate QDA params perfectly. If n too small and/or p too large LDA maybe preferred

\rightarrow Least amount of parameter estimations in almost all scenarios

Then it should not surprise that Gaussian NB often works best even though assumption of cond. feature independence is usually unrealistic!

Exercise Logistic Regression (Discriminant Approach)

hard label prediction $\hat{y} = 1 \Leftrightarrow \pi(x) \geq \alpha, \alpha \in (0,1)$

Proof that decision boundary is a linear hyperplane

Remember: Logistic regression estimates prob. $p(y=1|x, \theta) = \pi(x|\theta) = \frac{1}{1 + \exp(-\theta^T x)}$

$$\pi(x) = \frac{1}{1 + \exp(-\theta^T x)} = \alpha \Leftrightarrow 1 + \exp(-\theta^T x) = \frac{1}{\alpha} \Leftrightarrow \exp(-\theta^T x) = \frac{1}{\alpha} - 1 \Leftrightarrow -\theta^T x = \log\left(\frac{1}{\alpha} - 1\right) \hat{y} = 1 \text{ for all points on or above this plane}$$

linear comp. of features scalar b depending on α

For $\alpha=0.5$: $\theta^T x = -\log(2-1) = -\log(1) = 0$, so $\theta^T x \geq 0 \Leftrightarrow \hat{y} = 1$ "1" halfspace

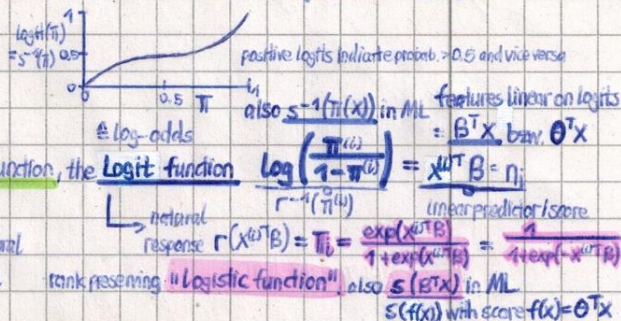
$\theta^T x < 0 \Leftrightarrow \hat{y} = 0$ "0" halfspace

Logistic Regression

Regression mit binärer Zielgrösse \Rightarrow Ein Parameter $\pi^{(i)} = P(y^{(i)} = 1) = E(y^{(i)})$ bzw. $P(y_i = 1 | x_i) = E(y_i | x_i)$
codiert als 1/0

Wir interessieren uns dafür, wie Einflussgrößen X mit Erwartungswert bzw. Wahrscheinlichkeit für das "1" Ereignis zusammenhängen

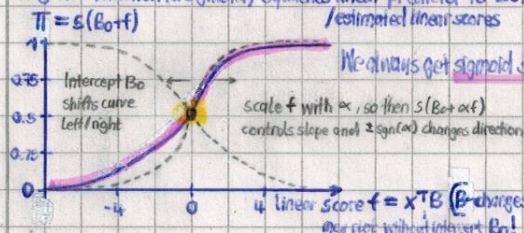
Bernoulli-Verteilung: $y^{(i)} | x^{(i)} \sim \mathcal{B}(\pi^{(i)})$ $i=1, \dots, n$ unabhängig



1 GLM Logistic Regression is simply Bernoulli-Regression with its natural Link Function, the Logit function $\log\left(\frac{\pi_i}{1-\pi_i}\right) = x_i^T \beta = \eta_i$

We could use other sigmoid types, and GOF would fit requirement, but logistic function is the most natural

Logistic function (a sigmoid) squashes linear predictor to $[0, 1]$ probability interval
 $\pi = g(\beta_0 + f)$ / estimated linear scores



with turning point at $(\text{score} = 0, \pi = 0.5)$, here maximal slope in region of highest doc uncertainty

$$s(0) = \frac{e}{2 \cdot \pi} \cdot \pi = \frac{1}{2} \quad \text{log-odds} = 0, \quad \log\left(\frac{0.5}{1-0.5}\right) = \log(1) = 0$$

natural decision boundary at $\pi = 0.5$ s.t. scores > 0 ($-\log(\text{odds}) > 0$) lead to class "1" and vice versa
however there are use cases where we don't base decision boundary on "inclusion" $\pi = 0.5$ (identical line)
boundary $= 0.5 \cdot \log(0.1/0.9) = \log(0.1/0.9) = 0$ equal treatment of $\log(\text{odds})$ misclassified observations
defined by x-axis
descent $B_0 - B_1$ slope also shape $f = x \cdot B$ as x-axis captures all linear combinations of features in one dimension

In 1D we could also plot feature on x-axis, then different μ_j change slope!
 ↳ Still sigmoidal shape with turning point at $(\tau = 0.5, \text{ but possibly } x_1 \neq 0)$

\Rightarrow have to group them to analyze, or $\epsilon_i = 0$ or 1 for each q_i .
Very small/large residuals at edges, medium in the middle
always $\epsilon_i = \hat{T}_i$ or $1 - \hat{T}_i$ and depend in size on curve location

GLM vs. classical LM: no residuals by model design & useless to analyze (follow no distribution)

Es besteht ein linearer Zusammenhang mit den log-Odds! Interpretation ist so wenig anders als bei klassischem linearem Modell für marginale c.p. Effekte von Δx

$$P(y=1|x) = \frac{\exp(x^T \beta)}{1 + \exp(x^T \beta)} \quad q = \frac{p}{1+p} \Rightarrow q(y=1|x) = \dots \exp(x^T \beta) \text{ estimated odds of class "1" for given feature vector } x$$

If we define $x^{+k} = x + \underset{\text{index } k}{(0, 0, \dots, 0, 1, 0, \dots, 0, 0)}$ then we get $q(y=1|x^{+k}) = \exp(B_k) q(y=1|x)$ $\log(q(y=1|x^{+k})) = \log(q(y=1|x)) + B_k$

multiplicative odds change when X_k increases by one unit additive log odds change
 in X for the log-odds of positive class (also odds ratio change by factor $\exp(\beta_k)$)

$\log(q(y=1|x)) = x^T \beta$ Logistic Regression assumes linear structure in x for the log-odds of positive class. (also odds ratio change by factor $\exp(\beta)$)

don't misinterpret $\exp(B^k)$ as relative risk!
 Only 2 risk factors \Rightarrow is very rare $\Rightarrow q() \approx p()$

MLE estimator in logistic function

log-likelihood

$$\hat{\beta}_{MLE} = \arg\max_{\beta} \ell(\beta) \text{ with } L(\beta) = \prod_{i=1}^n r(x_i^T \beta)^{y_i} (1 - r(x_i^T \beta))^{1-y_i} \text{ and } \ell(\beta) = \log L(\beta) = \sum y_i (\log(r(x_i^T \beta))) + (1-y_i) (\log(1-r(x_i^T \beta)))$$

$$\text{score-fcn. } s(B) = \frac{\partial \mathcal{L}}{\partial B} = \sum_{i=1}^n \left(w_i \frac{r(x^{(i)T}B)(1-r(x^{(i)T}B))x^{(i)}}{r(x^{(i)T}B)} + (1-w_i) \frac{-r(x^{(i)T}B)(1-r(x^{(i)T}B))x^{(i)}}{1-r(x^{(i)T}B)} \right) = \sum_{i=1}^n (w_i x^{(i)} - r(x^{(i)T}B)x^{(i)}) = \sum_{i=1}^n (w_i - r(x^{(i)T}B))x^{(i)}$$

$$r = \frac{\exp(x \cdot \beta)}{1 + \exp(x \cdot \beta)} = \frac{1}{1 + \exp(-x \cdot \beta)} = (1 + \exp(-x \cdot \beta))^{-1}$$

just taking the (-1) out of $x \cdot \beta$

without inner function chain rule: $r' = \frac{-(\exp(-x \cdot \beta) \cdot (-1))}{(1 + \exp(-x \cdot \beta))^2} = \frac{(\exp(-x \cdot \beta) - 1) \cdot 1}{(1 + \exp(-x \cdot \beta))^2} = \frac{r(1-r)}{1}$ here $\pi(1-\pi)$

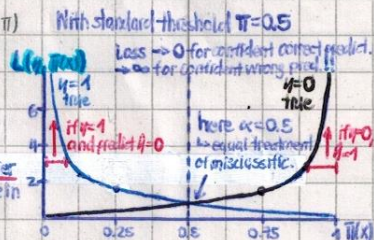
(-1) weil nach diff. $(\exp(-x \cdot \beta))$

⇒ For MLE parameter estimator: $S(\hat{\beta}_{MLE}) = \sum_{i=1}^n (y^{(i)} - r(x^{(i)T} \hat{\beta}_{MLE})) x^{(i)} = 0$

$$g_{\alpha\beta}^{\theta}$$
 is k-th diag. elem. of inverse Fisher matrix at point θ on $\mathcal{I}^{-1}(\theta)$

Wald confidence interval $\hat{\beta}_k \pm \hat{\sigma}_{\beta_k} z_{1-\alpha/2}$ with $\hat{\sigma}_{\beta_k} = \sqrt{f_{kk}}$

Log-loss penalizes heavier
the more confident we are in
a wrong prediction



For odds-ratios change-factor we get transformed CI $\exp[\hat{\beta}_k \pm \hat{\sigma}_{\hat{\beta}_k} z_{1-\alpha/2}]$

Free estimate $P(y=1|x, \theta)$ $\theta < 0 \Rightarrow$ hard label choice: $h(x^{(ii)} | \theta, \alpha) = \mathbb{I}_{[\alpha, 1]} \left(\frac{1}{1 + \exp(-\theta^T x^{(ii)})} \right)$ with $\alpha \in (0, 1)$, typically $\alpha = 0.5$

Hypothesis Space: $\mathcal{H} = \{\pi: X \rightarrow [0,1] \mid \pi(x) = \sigma(\theta^T x) = \frac{1}{1 + \exp(-\theta^T x)}, \theta \in \mathbb{R}^{n+1}\}$ with $f(x) = \theta^T x$ and $L(y, \pi) = -y f + \log(1 + \exp(f))$

Logistic Regression in Machine Learning: ERM requires suitable loss function to find optimal parameters $\hat{\theta} = \arg \min_{\theta} \text{Emp}(\theta) = \arg \min_{\theta} \sum_{i=1}^n L(y^{(i)}, f(x^{(i)}|\theta))$

$$P(y^{(1)}, \dots, y^{(n)}) = \prod_{i=1}^n P(y^{(i)} = y^{(i)}) \quad \text{with } y^{(i)} \in \{0, 1\} \text{ for negative/positive class}$$

Stori nfti: Likelihood $L(\theta) = \prod_{i=1}^n \pi(x^{(i)}|\theta) \prod_{i=1}^n (1 - \pi(x^{(i)}|\theta)) = \prod_{i=1}^n \pi(x^{(i)}|\theta) (1 - \pi(x^{(i)}|\theta))$

Now log-to-turn nasty product into nice sum $\ell(\theta) = \log(L(\theta)) = \sum_{i=1}^n (y_i^{(i)} \log(\pi(x^{(i)} | \theta)) + (1 - y_i^{(i)}) \log(1 - \pi(x^{(i)} | \theta)))$ maximize this to find $\hat{\theta}$

④ maximizes Likelihood & minimizes Rmse of model

... to turn this into minimization problem: $-l(\theta) = \sum_{i=1}^n -y^{(i)} \log(\pi(x^{(i)}|\theta)) - (1-y^{(i)}) \log(1-\pi(x^{(i)}|\theta))$ $\log(\pi) = \log\left(\frac{\exp(-f)}{1+\exp(-f)}\right) = -f \log(1+\exp(-f)) - \log\left(\frac{\exp(-f)}{1+\exp(-f)}\right)$

also sometimes called cross-entropy loss, works toward Log Reg. for all types of binary classification models with Discriminant ERM approach $L(y, \hat{y}) = -(1-y)f + \log(1 + \exp(f))$ $\frac{\log(1 + \exp(f))}{\log(1 + \exp(f))} = f + \frac{1}{\log(1 + \exp(f))}$
 \Rightarrow Bernoulli/Binomial/Log-Loss $L(y, \hat{y}) = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$ with logit link $\hat{y} = \frac{1}{1 + \exp(-f(x))} = \sigma(f(x)) \Leftrightarrow$ score-based $L(y, f) = -yf + \log(1 + \exp(f))$

Convex ✓ Logistic Regression has unique solution under simple regularity conditions, so design matrix X must have full rank and det. not perfectly linearly dependent (also semi-def. linear) ✓ Numerical optimization with 2nd order f.t.e. a Newton-Raphson if linear is not too computationally unbounded, usually 20-30 it's fine ✓

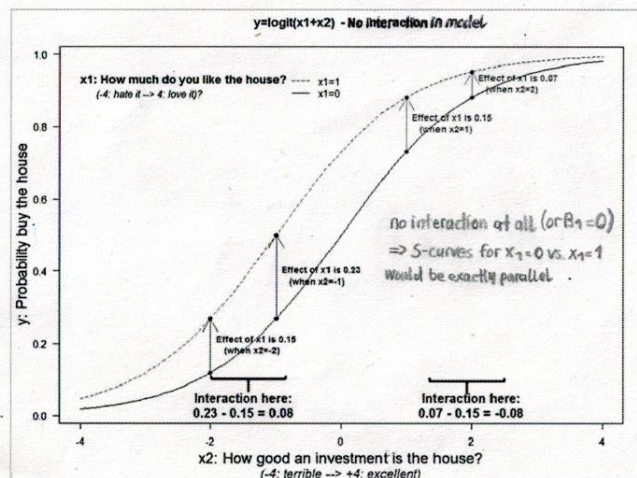
Interpretation von logistischer Regression

Die folgenden Aussagen gelten für eine Einflussgröße x_k bei Festhalten aller anderen Einflussgrößen (c.p.)

- Erhöht sich x_k um eine Einheit, so ändert sich die logarithmierte Chance von Y um β_k .
- Wenn x_k um eine Einheit steigt, so ändert sich die Chance von Y um den Faktor $\exp(\beta_k)$.
- Das Odds Ratio (Chancenverhältnis) zwischen Y bei x_k und Y bei $x_k + 1$ ist $\exp(\beta_k)$.

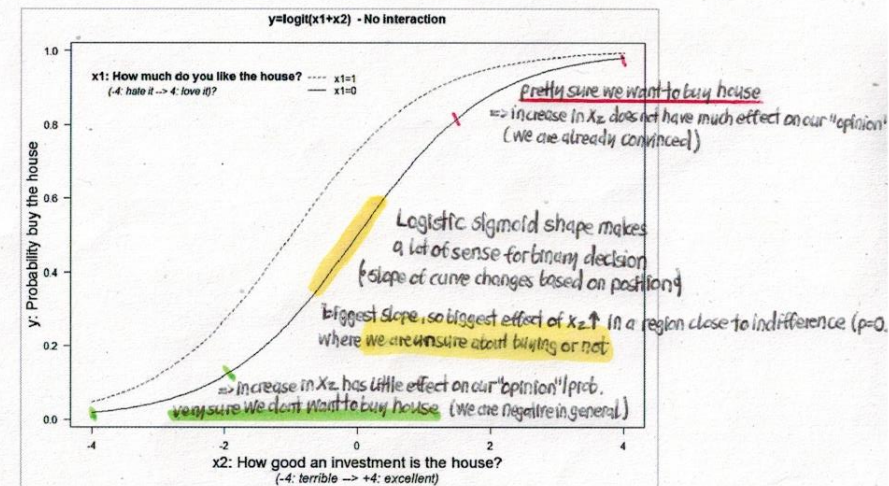
W'keit	0.01	0.05	0.1	0.3	0.4	0.5	0.6	0.7	0.9	0.95	0.99
Odds	1/99	1/19	1/9	3/7	2/3	1	1.5	7/3	9	19	99
Log odds	-4.6	-2.9	-2.2	-0.85	-0.41	0	0.41	0.85	2.2	2.9	4.6

Logistische Regression: Interpretation von Interaktionen



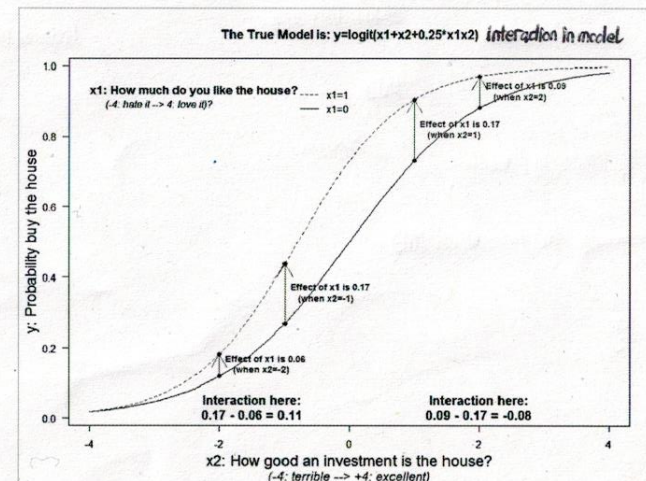
Source: <http://datacolada.org/57>

Logistische Regression: Logistische Link-Funktion



Source: <http://datacolada.org/57>

Logistische Regression: Interpretation von Interaktionen



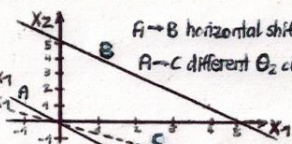
Source: <http://datacolada.org/57>

Linear classifiers like logistic regression learn a decision boundary that takes the form of a (linear) hyperplane.

Hyperplanes can be defined by equations $\theta^T x = b \Leftrightarrow \theta^T x - b = 0$ with coefficients θ and scalar $b \in \mathbb{R}$.

In order to see that such expressions actually describe hyperplanes, consider $\theta^T x = \theta_0 + \theta_1 x_1 + \theta_2 x_2 = 0$. Sketch the hyperplanes given by the following coefficients and explain the difference between the parameterizations:

- A • $\theta_0 = 0, \theta_1 = \theta_2 = 1$ $x_1 + x_2 = 0, x_2 = -x_1$
 B • $\theta_0 = -5, \theta_1 = \theta_2 = 1$ $-5 + x_1 + x_2 = 0, x_2 = +5 - x_1$
 C • $\theta_0 = 0, \theta_1 = 1, \theta_2 = 2$ $x_1 + 2x_2 = 0, x_2 = -\frac{1}{2}x_1$



Exercise 3: Decision Boundaries & Thresholds in Logistic Regression

Learning goals

- 1) Understand that logistic regression finds a linear decision boundary
- 2) Get a feeling for how parameterization changes predicted probabilities

In logistic regression (binary case), we estimate the probability $p(y=1|x, \theta) = \pi(x|\theta)$. In order to decide about the class of an observation, we set $\hat{y} = 1$ iff $\pi(x|\theta) \geq \alpha$ for some $\alpha \in (0, 1)$.

Show that the decision boundary of the logistic classifier is a (linear) hyperplane.

Hint

Derive the value of $\theta^T x$ (depending on α) starting from which you predict $\hat{y} = 1$ rather than $\hat{y} = 0$.

Below you see the logistic function for a binary classification problem with two input features for different values $\theta^T = (\theta_1, \theta_2)^T$ (plots 1-3) as well as α (plot 4). What can you deduce for the values of θ_1, θ_2 , and α ? What are the implications for classification in the different scenarios?

Derive the equation for the decision boundary hyperplane if we choose $\alpha = 0.5$.

Explain when it might be sensible to set α to 0.5:

$\alpha = 0.5$: losses of misclassified observations are treated equally \uparrow Log-loss $L(y, \pi) = -y \log(\pi) - (1-y) \log(1-\pi)$
 i.e. $L(y=0|\pi=1) = -\log(\pi)$ and $L(y=1|\pi=0) = -\log(1-\pi)$ are identical at boundary $\pi(x|\theta) = 0.5$!

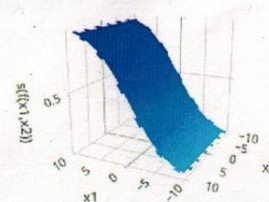
$\alpha \gg 0.5$: (choose this if we need be careful with class "1" predictions

For example when decision triggers a costly/dangerous therapy, so we want to minimize False-Positives

Vice versa: $\alpha \ll 0.5$ if we want to heavily emphasize penalty on False-Negatives, e.g. airport weapon security

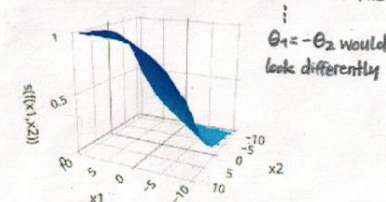
Plot (1)

- hypersurface is symmetric parallel across x_2 -axis
- $\Rightarrow \theta_2 = 0, x_2$ has zero effect
- $\theta_1 \approx 0.5$ by visual cue



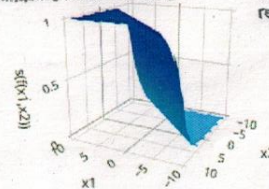
Plot (2)

- both feature dimensions affect logistic function, to equal degree $\Rightarrow \theta_1 = \theta_2$; hypersurface not tilted in either direction
- We also see that $\theta_1, \theta_2 > 0$ as $s(\cdot)$ \uparrow with $x_1 \uparrow, x_2 \uparrow$



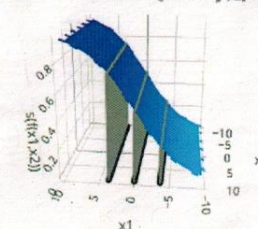
Plot (3)

- Very similar to (2) with minor differences:
- steeper logistic fit \Rightarrow larger abs. θ_1, θ_2 values
- sharper separation between class predictions, less in the more confidence in our decisions $\pi(x) = 0.5$ region



Plot (4)

- This is Plot (1) with different α -thresholds visualised
- midline at $x_1 = 0$ corresponds to $\alpha = 0.5$
- Other hyperplanes \neq higher (left) and lower (right) thresholds



misleading slices, just vertical

- it appears as if x_1 is const. at thresholds in general
- actually its just a geometric representation which looks like this here because of $\theta_2 = 0$