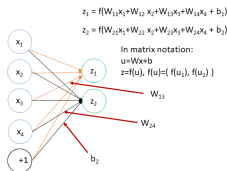


Chapter 3: Introduction to Machine Learning – Linear and Logistic Regression



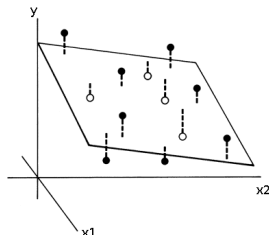
Christian Heumann

(largely based on slides from Nina Poerner, Benjamin Roth, Marina Speranskaya)

CIS LMU München, Department of Statistics LMU München

November 2020

Linear Regression



- In linear regression, the prediction \hat{y} of a label y is a linear function:

$$\hat{y} = \mathbf{w}^T \mathbf{x} = \sum_{j=1}^n w_j x_j$$

- Note, that in this notation, the parameter vector $\mathbf{w} \in \mathbb{R}^n$ stands already for the *estimated* parameter vector.
- Weight w_j decides if value of feature x_j increases or decreases prediction \hat{y} .

Matrix notation I

- m samples
- Prediction for one sample:

$$\hat{y}_i = \mathbf{w}^T \mathbf{x}_i$$

- Error for one sample (residual):

$$e_i = (y_i - \hat{y}_i)$$

- Squared error for one sample:

$$e_i^2 = (y_i - \hat{y}_i)^2 = (\hat{y}_i - y_i)^2$$

Matrix notation II

- Matrix notation:

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \quad \hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_m^T \end{bmatrix} \quad \hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$$

- \mathbf{X} in detail: typically, first column contains all **1** for the intercept (bias, shift) term.

$$\mathbf{X} = \begin{bmatrix} 1 & x_{12} & x_{13} & \dots & x_{1n} \\ 1 & x_{22} & x_{23} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m2} & x_{m3} & \dots & x_{mn} \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

- Estimate parameters using $\mathbf{X}^{(train)}$ and $\mathbf{y}^{(train)}$.
- Make high-level decisions (which features...) using $\mathbf{X}^{(dev)}$ and $\mathbf{y}^{(dev)}$.
- Evaluate resulting model using $\mathbf{X}^{(test)}$ and $\mathbf{y}^{(test)}$.

Simple Example: Housing Prices

- Predict property prices (in 1K Euros) from just one feature: Square feet of property.

$$\mathbf{X} = \begin{bmatrix} 1 & 450 \\ 1 & 900 \\ 1 & 1350 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 730 \\ 1300 \\ 1700 \end{bmatrix}$$

- Prediction is:

$$\hat{\mathbf{y}} = \begin{bmatrix} w_1 + 450w_2 \\ w_1 + 900w_2 \\ w_1 + 1350w_2 \end{bmatrix} = \begin{bmatrix} 1 & 450 \\ 1 & 900 \\ 1 & 1350 \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \mathbf{X}\mathbf{w}$$

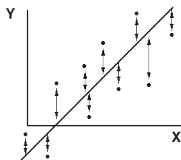
- \mathbf{w}_1 will contain costs incurred in any property acquisition
- \mathbf{w}_2 will contain remaining average price per square feet.
- Optimal parameters are for the above case:

$$\mathbf{w} = \begin{bmatrix} 273.3 \\ 1.08 \end{bmatrix} \quad \hat{\mathbf{y}} = \begin{bmatrix} 759.1 \\ 1245.1 \\ 1731.1 \end{bmatrix}$$

Linear Regression: Mean Squared Error

- Mean squared error of training (or test) data set is the sum of squared differences between the predictions and labels of all m instances.

$$MSE := \frac{1}{m} \sum_{i=1}^m e_i^2 = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$



- In matrix notation:

$$\begin{aligned} MSE &:= \frac{1}{m} \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2 \\ &= \frac{1}{m} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 \end{aligned}$$

Learning: Improving on MSE

- Gradient of a function f : Vector whose components are the n partial derivatives of f wrt to the parameters, here \mathbf{w} .

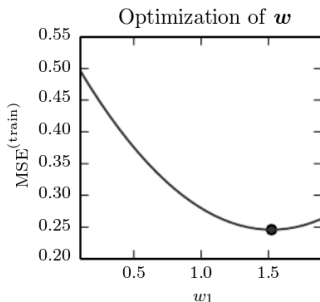
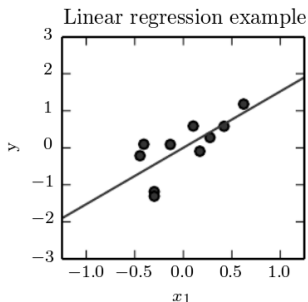
$$\nabla_{\mathbf{w}} f(\mathbf{w}) = \begin{bmatrix} \frac{\partial f(\mathbf{w})}{\partial w_1} \\ \frac{\partial f(\mathbf{w})}{\partial w_2} \\ \vdots \\ \frac{\partial f(\mathbf{w})}{\partial w_n} \end{bmatrix}$$

- View MSE as a function $f(\mathbf{w})$ of \mathbf{w}
- Minimum is where gradient is $\mathbf{0}$:

$$\nabla_{\mathbf{w}} MSE \stackrel{!}{=} \mathbf{0}$$

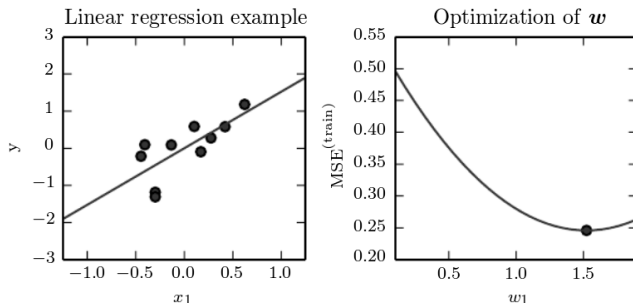
Learning: Improving on MSE

- View MSE as a function of \mathbf{w}



- Minimum is where gradient $\nabla_{\mathbf{w}} MSE = \mathbf{0}$.
- Why minimum and not maximum or saddle point?
 - ▶ Because it is a quadratic function...
 - ▶ Check convexity for 1 dimensional function: Second derivative > 0 .
 - ▶ Check for vector valued function: Hessian is positive-semidefinite.

Second Derivative Test



Second derivative of Mean Squared Error for Linear model with only one feature:

$$\frac{d^2}{dw^2} \sum_{i=1}^m (x^{(i)}w - y^{(i)})^2 = \frac{d^2}{dw^2} \sum_{i=1}^m (x^{(i)2}w^2 - 2x^{(i)}w + y^{(i)2}) = 2 \sum_{i=1}^m x^{(i)2} > 0$$

Solving for \mathbf{w}

- We now know that minimum is where gradient is $\mathbf{0}$.

$$\nabla_{\mathbf{w}} MSE = \mathbf{0}$$

$$\Rightarrow \nabla_{\mathbf{w}} \frac{1}{m} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = \mathbf{0}$$

- Solve for \mathbf{w} :

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

(Normal Equation)

- The inverse $(\mathbf{X}^T \mathbf{X})^{-1}$ exists, if X has full column rank (i.e. rank n).

Deriving the Normal Equation

- Function to minimize:

$$\begin{aligned} & ||\mathbf{X}\mathbf{w} - \mathbf{y}||_2^2 \\ &= (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \\ &= \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y} \\ &= \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y} \end{aligned}$$

- Take the gradient¹ w.r.t. \mathbf{w} and set equal to $\mathbf{0}$:

$$\begin{aligned} 2\mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{X}^T \mathbf{y} &= \mathbf{0} \\ \Rightarrow \mathbf{X}^T \mathbf{X} \mathbf{w} &= \mathbf{X}^T \mathbf{y} \\ \Rightarrow (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \mathbf{w} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

¹[Matrix Cookbook. Petersen and Pedersen, 2012]:

$$\nabla_{\mathbf{w}} \mathbf{w}^T \mathbf{a} = \mathbf{a}$$

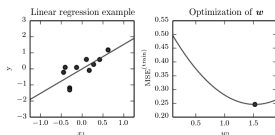
$$\nabla_{\mathbf{w}} \mathbf{w}^T \mathbf{B} \mathbf{w} = 2\mathbf{B} \mathbf{w} \text{ for symmetric } \mathbf{B}$$

Practical Remarks

What if $\mathbf{X}^T \mathbf{X}$ does not have an inverse?

- This can happen if there are infinitely many solutions:
 - ▶ one feature is the exact multiple of another
 - ▶ there are more features than training examples
 - ▶ \Rightarrow a Moore-Penrose pseudoinverse picks solution with smallest Euclidean norm.
 - ▶ \Rightarrow adding a regularization term makes the system of equations non-singular (uniquely solvable).
 - ★ Normal Equation becomes: $\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$
 - ★ Also called *Ridge Regression*.
- In general: **Avoid computing the matrix inverse when implementing least squares (slow/unstable)!**
- There are usually special routines for solving linear least squares and ridge regression.

Second Derivative Test for Multi-Dimensional Case



- Second derivative test for multi-dimensional $\mathbf{w} \in \mathbb{R}^n$
Is Hessian positive-semidefinite? ($\mathbf{z}^T \mathbb{H} \mathbf{z} > 0$ for $\mathbf{z} \neq \mathbf{0}$)
 \Leftrightarrow Function is convex.
- Hessian: Matrix of second-order partial derivatives.

$$\mathbb{H}_{j,k} = \frac{\partial^2 f(\mathbf{w})}{\partial w_j \partial w_k}$$

- Matrix algebra² gives:

$$\mathbb{H} = 2\mathbf{X}^T \mathbf{X} = 2 \sum_{i=1}^m \mathbf{x}^{(i)} \mathbf{x}^{(i)T}; \quad \mathbf{z}^T \mathbb{H} \mathbf{z} = 2 \sum_{i=1}^m (\mathbf{z}^T \mathbf{x}^{(i)})^2 > 0$$

²[Matrix Cookbook. Petersen and Pedersen, 2012]

Linear Regression: Summary

- Linear regression models simple linear relationships between \mathbf{X} and \mathbf{y}
- The Mean squared error is a quadratic function of the parameter vector \mathbf{w} , and has a unique minimum.
- Normal equations: Find the minimum by setting the gradient to zero and solving for \mathbf{w} .
- Linear algebra packages have special routines for solving least squares linear regression.

Maximum Likelihood Estimation

- Machine learning models are often more interpretable if they are stated in a probabilistic way.
- Performance measure: What is the probability of the training data given the model parameters? Works only well for discrete data! More general: use densities instead of probabilities
- Likelihood: Probability density of data as a function of model parameters. Generally, likelihood is a function proportional to the density.
- \Rightarrow Maximum Likelihood Estimation
- Many models can be formulated in a probabilistic way!

Probability density of a data set

- Data:

- ▶ Set of m examples $\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$
- ▶ Sometimes written as design matrix:

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}^{(1)T} \\ \vdots \\ \mathbf{x}^{(m)T} \end{bmatrix}$$

- Probability density of a data set \mathbf{X} , parametrized by θ :

$$p_{model}(\mathbf{X}; \theta)$$

Probability density of a data set

- Data points are assumed to be (stochastically) independent (and sometimes identically distributed) random variables (i.d. or i.i.d.)
 - ▶ Assumption made by many ML models.
 - ▶ Identically distributed: Examples come from same distribution.
 - ▶ Independent: Value of one example doesn't influence other example.
 - ▶ \Rightarrow Probability density of the data set is the product of example probability densities.

$$p_{model}(\mathbf{X}; \theta) = \prod_{i=1}^m p_{model}(\mathbf{x}^{(i)}; \theta)$$

Maximum Likelihood Estimation

- Likelihood: Probability density of data viewed as function of parameters θ
- (Negative) Log-Likelihood (NLL):
 - ▶ Logarithm is monotonically increasing
 - ★ Maximum of function stays the same
 - ★ Easier to do arithmetic with (sums vs. products)
 - ▶ Optimization is often formulated as minimization \Rightarrow take negative of function.
- Maximum likelihood estimator for θ :

$$\begin{aligned}\theta_{ML} &= \operatorname{argmax}_{\theta} p_{\text{model}}(\mathbf{X}; \theta) \\ &= \operatorname{argmax}_{\theta} \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \theta) \\ &= \operatorname{argmax}_{\theta} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \theta)\end{aligned}$$

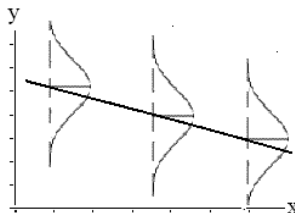
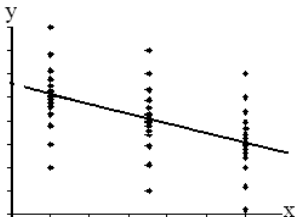
Conditional Log-Likelihood

- Log-likelihood can be stated for supervised and unsupervised tasks.
- Unsupervised learning (e.g. density estimation).
 - ▶ Task: model $p_{model}(\mathbf{X}; \theta)$ (as before)
 - ▶ $\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$
- Supervised learning (Predictive modeling):
 - ▶ Task: model $p_{model}(\mathbf{y}|\mathbf{X}; \theta)$
 - ▶ $\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$, $\mathbf{y} = \{y^{(1)}, \dots, y^{(m)}\}$
- Maximum likelihood estimation for the supervised case (independent examples):

$$\begin{aligned}\theta_{ML} &= \operatorname{argmax}_{\theta} P(\mathbf{y}|\mathbf{X}; \theta) \\ &= \operatorname{argmax}_{\theta} \sum_{i=1}^m \log P(y^{(i)}|\mathbf{x}^{(i)}; \theta)\end{aligned}$$

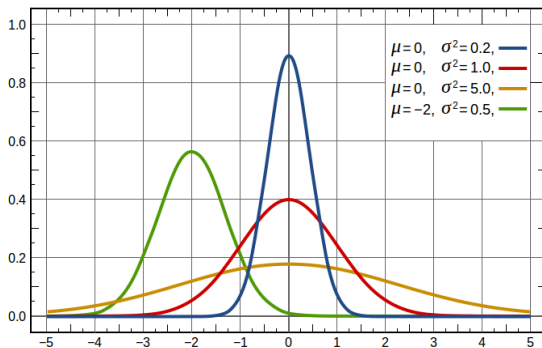
Linear Regression as Maximum Likelihood

- Instead of predicting one value \hat{y} for an input x , model probability distribution $p(y|x)$.
- For the same value of x , different values of y may occur (with different probability).



Gaussian Distribution

- Gaussian distribution: $N(y|\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(y-\mu)^2}{2\sigma^2}\right]$
 - ▶ Quadratic function as negative exponent, scaled by variance
 - ▶ Normalization factor $\frac{1}{\sigma\sqrt{2\pi}}$



Linear Regression as Maximum Likelihood

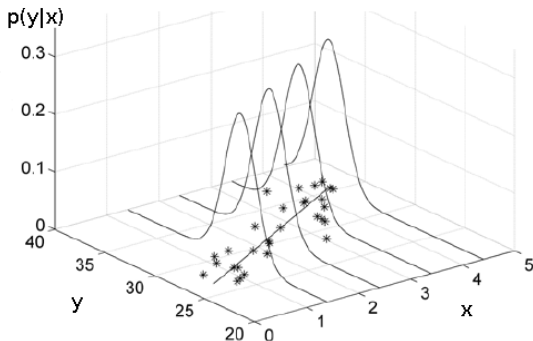
- Assume label y is distributed by a Gaussian, depending on features \mathbf{x}

$$p(y|\mathbf{x}) = N(y|\mu, \sigma^2)$$

where the mean is determined by the linear transformation

$$\mu = \boldsymbol{\theta}^T \mathbf{x}$$

and σ is a constant.



Linear Regression as Maximum Likelihood

- Gaussian distribution: $N(y|\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(y-\mu)^2}{2\sigma^2}\right]$
 - ▶ Taking the log makes it a quadratic function!
- Conditional negative log-likelihood:

$$\begin{aligned} & -\log P(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta}) \\ &= -\sum_{i=1}^m \log p(y^{(i)}|\mathbf{x}^{(i)}; \boldsymbol{\theta}) \\ &= m \log \sigma + \frac{m}{2} \log(2\pi) + \sum_{i=1}^m \frac{(y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)})^2}{2\sigma^2} \\ &= \text{const} + \frac{1}{2\sigma^2} \sum_{i=1}^m (y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)})^2 \end{aligned}$$

- What is the optimal value for $\boldsymbol{\theta}$?

Linear Regression as Maximum Likelihood

- Conditional negative log-likelihood:

$$NLL(\boldsymbol{\theta}) = \text{const} + \frac{1}{2\sigma^2} \sum_{i=1}^m (y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)})^2$$

- Compare to previous result:

$$MSE(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^n (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2$$

- Minimizing NLL under these assumptions is equivalent to minimizing MSE!
- The result of the minimization of the NLL is the MLE $\hat{\boldsymbol{\theta}}$.

A Caveat: Maximum Likelihood is not *Bayesian*

- MLE is simple: Identify distribution and parameters, then maximize!
- MLE accounts for (some) uncertainty: Instead of predicting a single value \hat{y} , treat y as a random variable.
- However, what about the uncertainty about our parameters θ ?
- Shouldn't θ be treated as a random variable, too?
(Instead of a point estimate $\hat{\theta}$?)

A Caveat: Maximum Likelihood is not *Bayesian*

- Uncertainty about θ influenced by:
 - ▶ Our assumptions what reasonable values for θ look like.
 - ▶ The amount of training data.
 - ▶ The properties (variance ...) of the training data.
- *Bayesian inference* is all about modelling randomness of θ in a sound way.
- Why is it called *Bayesian*?

$$p(\theta|\mathbf{X}) = \frac{p(\mathbf{X}|\theta)p(\theta)}{p(\mathbf{X})} = \frac{p(\mathbf{X}|\theta)p(\theta)}{\int p(\mathbf{X}|\theta)p(\theta)d\theta}$$

- The Bayesian approach is theoretically more appealing than working with point estimates.
- In practical terms: Sometimes going the Bayesian way pays off, sometimes it doesn't.

Maximum Likelihood: Summary

- Many machine learning problems can be stated in a probabilistic way.
- Mean squared error linear regression can be stated as a probabilistic model that allows for Gaussian random noise around the predicted value \hat{y} .
- A straightforward optimization is to maximize the likelihood of the training data.
- Maximum likelihood is not Bayesian, and may give undesirable results (e.g. if there is only little training data).
- In practice, MLE and point estimates are often used to solve machine learning problems.

Linear regression and mean squared error of prediction

- The linear model is given by (Y is a random variable, y its realization)

$$Y = \boldsymbol{\theta}^T \mathbf{x} + \varepsilon$$

with $\varepsilon \sim N(0, \sigma^2)$. Thus the *expectation* of Y is $E(Y) = \boldsymbol{\theta}^T \mathbf{x}$.

- Using the estimated parameter (Least Squares or MLE) $\hat{\boldsymbol{\theta}}$, the predicted value \hat{y} is given by

$$\hat{y} = \hat{\boldsymbol{\theta}}^T \mathbf{x}$$

which is also the *estimated* expectation of y , $\widehat{E(Y)}$.

- The theoretical quantity, which is estimated by the introduced MSE, is the so-called MSEP (*mean squared error of prediction*), which is also an expectation, based on a *new* observation Y_{m+1} (we ignore the index in the following) and its prediction $\hat{y} = \hat{\boldsymbol{\theta}}^T \mathbf{x}_{m+1}$:

$$MSEP = E \left\{ (Y - \hat{y})^2 \right\} = E \left\{ (Y - E(Y) + E(Y) - \hat{y})^2 \right\} = \text{Var}(Y) + MSE(\hat{y}) .$$

Note, that some terms are zero, because of stochastic independence of Y_{m+1} and \hat{y} . Note, that $MSEP > \text{Var}(Y) = \sigma^2$.

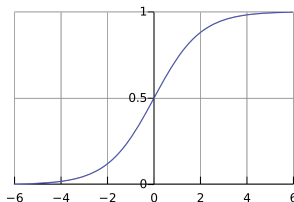
From Regression to Classification

- So far, linear regression:
 - ▶ A simple linear model.
 - ▶ Probabilistic interpretation.
 - ▶ Find optimal parameters using Maximum Likelihood Estimation.
- Can we do something similar for classification?
- \Rightarrow Logistic Regression (*... it's not actually used for regression ...*)

Logistic Regression

- Binary logistic model:
Estimate the probability of a binary response $y \in \{0, 1\}$ based on features \mathbf{x} .
- Logistic Regression is a *Generalized Linear Model*:
Linear model is related to the response variable via a **link function**.

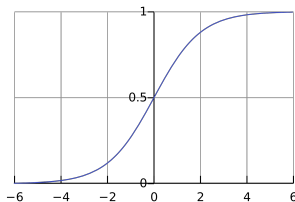
$$p(Y = 1|\mathbf{x}; \boldsymbol{\theta}) = f(\boldsymbol{\theta}^T \mathbf{x})$$



(Note: Y denotes a random variable, whereas y , $y^{(i)}$, $0, 1$ denote values that the random variable can take on. If the random variable is obvious from the context, it may be omitted.)

Logistic Regression

- Recall linear regression: $p(y|\mathbf{x}; \boldsymbol{\theta}) = N(y; \boldsymbol{\theta}^T \mathbf{x}, \sigma^2 \mathbf{I})$
 - Predicts $y \in \mathbb{R}$
- Classification: Outcome (per example) 0 or 1
 - Logistic sigmoid: $\sigma(z) = \frac{1}{1+e^{-z}}$



- Logistic Regression: Linear function + logistic sigmoid

$$p(Y = 1|\mathbf{x}; \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^T \mathbf{x})$$

Binary Logistic Regression

Probability of different outcomes (for one example):

- Probability of positive outcome:

$$p(Y = 1|\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}}}$$

- Probability of negative outcome:

$$p(Y = 0|\mathbf{x}; \boldsymbol{\theta}) = 1 - p(Y = 1|\mathbf{x}; \boldsymbol{\theta})$$

Probability of a Training Example

- Probability for actual label $y^{(i)}$ given features $\mathbf{x}^{(i)}$
- Can be written for both labels (0 and 1) without case distinction
- Label exponentiation trick: use $x^0 = 1$

$$\begin{aligned} & p(Y = y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) \\ &= \begin{cases} p(Y = 1 | \mathbf{x}^{(i)}; \boldsymbol{\theta}) & \text{if } y^{(i)} = 1 \\ p(Y = 0 | \mathbf{x}^{(i)}; \boldsymbol{\theta}) & \text{if } y^{(i)} = 0 \end{cases} \\ &= \begin{cases} p(Y = 1 | \mathbf{x}^{(i)}; \boldsymbol{\theta})^1 p(Y = 0 | \mathbf{x}^{(i)}; \boldsymbol{\theta})^0 & \text{if } y^{(i)} = 1 \\ p(Y = 1 | \mathbf{x}^{(i)}; \boldsymbol{\theta})^0 p(Y = 0 | \mathbf{x}^{(i)}; \boldsymbol{\theta})^1 & \text{if } y^{(i)} = 0 \end{cases} \\ &= p(Y = 1 | \mathbf{x}^{(i)}; \boldsymbol{\theta})^{y^{(i)}} p(Y = 0 | \mathbf{x}^{(i)}; \boldsymbol{\theta})^{1-y^{(i)}} \end{aligned}$$

Binary Logistic Regression

- Conditional Negative Log-Likelihood (NLL):

$$\begin{aligned}NLL(\boldsymbol{\theta}) &= -\log p(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta}) \\&= -\log \prod_{i=1}^m p(Y = y^{(i)}|\mathbf{x}^{(i)}; \boldsymbol{\theta}) \\&= -\log \prod_{i=1}^m p(Y = 1|\mathbf{x}^{(i)}; \boldsymbol{\theta})^{y^{(i)}} (1 - p(Y = 1|\mathbf{x}^{(i)}; \boldsymbol{\theta}))^{1-y^{(i)}} \\&= -\sum_{i=1}^m y^{(i)} \log \sigma(\boldsymbol{\theta}^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log(1 - \sigma(\boldsymbol{\theta}^T \mathbf{x}^{(i)}))\end{aligned}$$

- No closed form solution for minimum
- Use numerical / iterative methods.
- LBFGS, Gradient descent ...

Logistic Regression

- Logistic regression: Logistic sigmoid function applied to a weighted linear combination of feature values.
- To be interpreted as the probability that the label for a specific example equals 1.
- Applying the model on test data: Predict $y^{(i)} = 1$ if

$$p(Y = 1|\mathbf{x}^{(i)}; \boldsymbol{\theta}) > 0.5$$

- No closed form solution for maximizing NLL, iterative methods necessary.