

10-601

Machine Learning

Logistic regression

Back to classification

1. Instance based classifiers

- Use observation directly (no models)
- e.g. K nearest neighbors

2. Generative:

- build a generative statistical model
- e.g., Bayesian networks

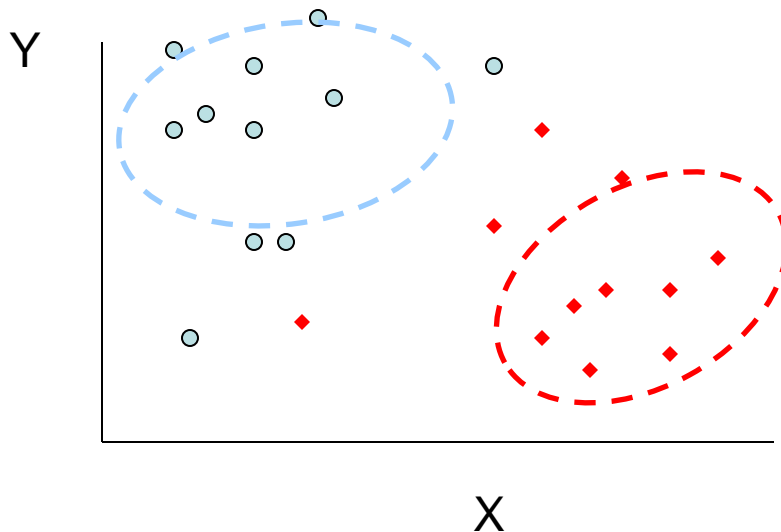
3. Discriminative

- directly estimate a decision rule/boundary
- e.g., decision tree

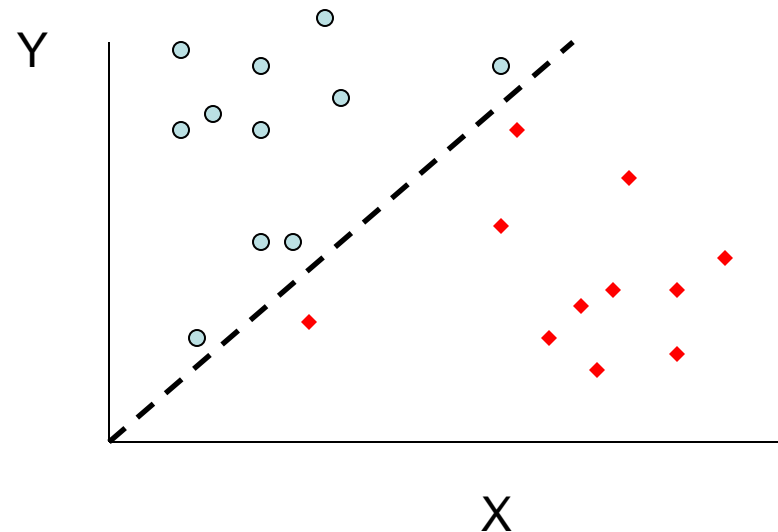
Generative vs. discriminative classifiers

- When using generative classifiers we relied on all points to learn the generative model
- When using discriminative classifiers we mainly care about the boundary

Generative model



Discriminative model



Regression for classification

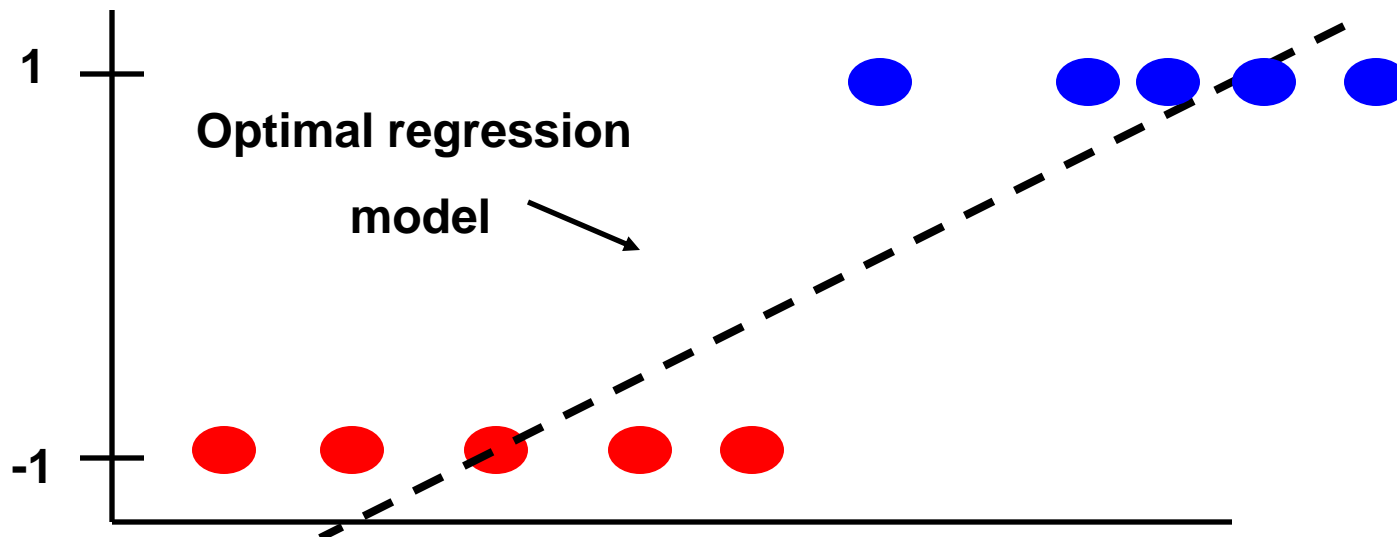
- In some cases we can use linear regression for determining the appropriate boundary.
- However, since the output is usually binary or discrete there are more efficient regression methods
- Recall that for classification we are interested in the conditional probability $p(y | x ; \theta)$ where θ are the parameters of our model
- When using regression θ represents the values of our regression coefficients (w).

Regression for classification

- Assume we would like to use linear regression to learn the parameters for $p(y | x ; \theta)$
- Problems?

$$\mathbf{w}^T \mathbf{x} \geq 0 \Rightarrow \text{classify as } 1$$

$$\mathbf{w}^T \mathbf{x} < 0 \Rightarrow \text{classify as } -1$$



The sigmoid function

$$p(y | x; \theta)$$

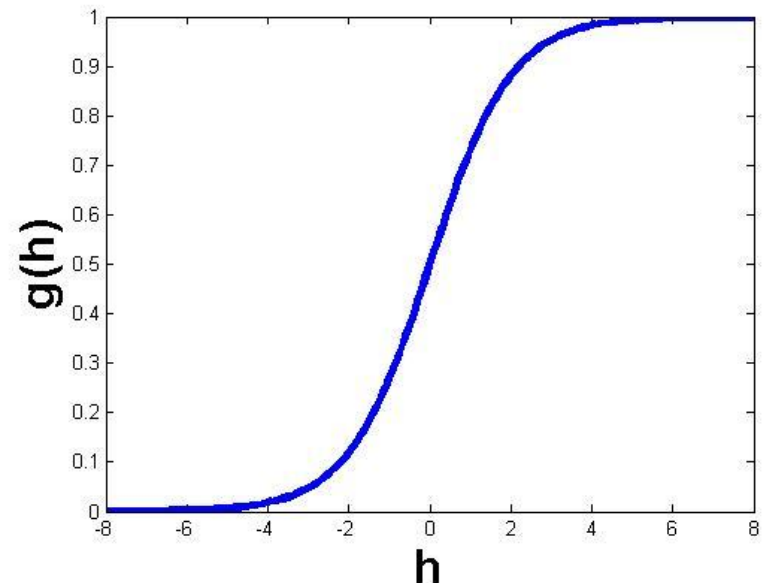
- To classify using regression models we replace the linear function with the sigmoid function:

Always between 0 and 1 \longrightarrow $g(h) = \frac{1}{1 + e^{-h}}$

- Using the sigmoid we set (for binary classification problems)

$$p(y = 0 | x; \theta) = g(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x}}}$$

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



The sigmoid function

$$p(y | x; \theta)$$

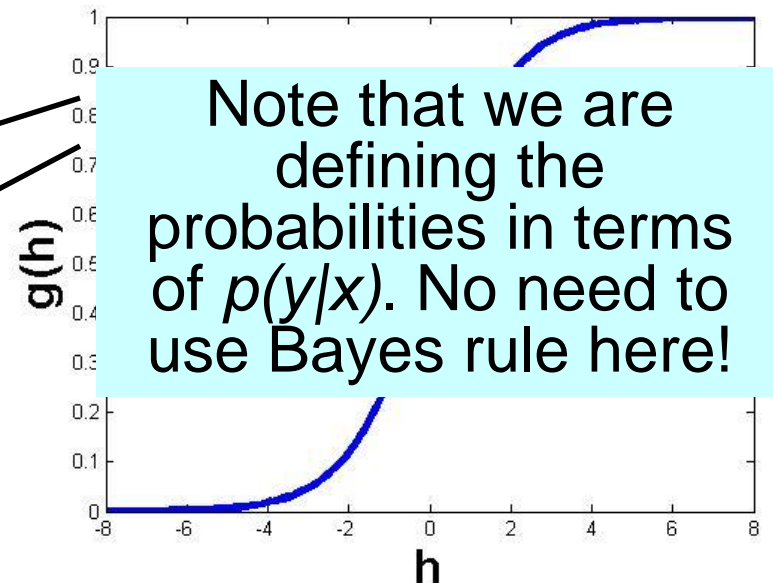
- To classify using regression models we replace the linear function with the sigmoid function:

$$g(h) = \frac{1}{1 + e^{-h}}$$

- Using the sigmoid we set (for binary classification problems)

$$p(y = 0 | x; \theta) = g(w^T x) = \frac{1}{1 + e^{w^T x}}$$

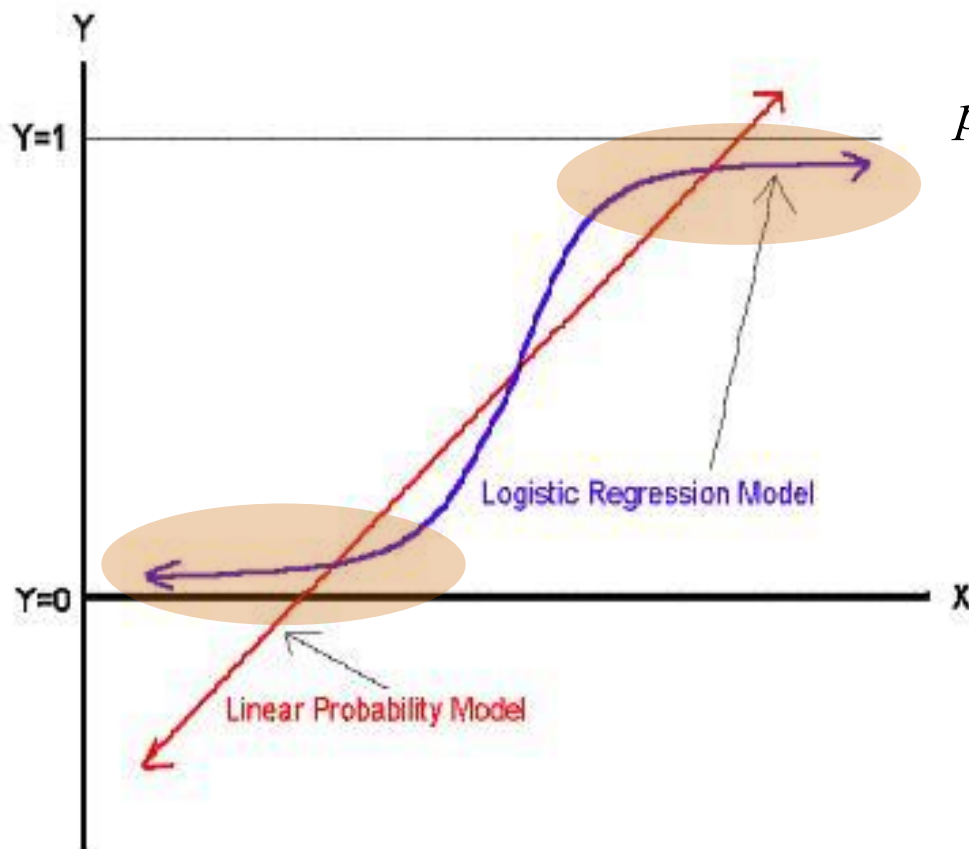
$$p(y = 1 | x; \theta) = 1 - g(w^T x) = \frac{e^{w^T x}}{1 + e^{w^T x}}$$



Logistic regression vs. Linear regression

$$p(y = 0 | x; \theta) = g(w^T x) = \frac{1}{1 + e^{w^T x}}$$

$$p(y = 1 | x; \theta) = 1 - g(w^T x) = \frac{e^{w^T x}}{1 + e^{w^T x}}$$



Determining parameters for logistic regression problems

- So how do we find the parameters?
- Similar to other regression problems we look for the MLE for w
- The likelihood of the data given the model is:

$$p(y = 0 | x; \theta) = g(x; w) = \frac{1}{1 + e^{w^T x}}$$

$$p(y = 1 | x; \theta) = 1 - g(x; w) = \frac{e^{w^T x}}{1 + e^{w^T x}}$$

$$L(y | x; w) = \prod_i (1 - g(x^i; w))^{y^i} g(x^i; w)^{(1-y^i)}$$

Solving logistic regression problems

$$g(x; w) = \frac{1}{1 + e^{w^T x}}$$

$$1 - g(x; w) = \frac{e^{w^T x}}{1 + e^{w^T x}}$$

- The likelihood of the data is: $L(y | x; w) = \prod_i (1 - g(x^i; w))^{y^i} g(x^i; w)^{(1 - y^i)}$
- Taking the log we get:

$$\begin{aligned} LL(y | x; w) &= \sum_{i=1}^N y^i \ln(1 - g(x^i; w)) + (1 - y^i) \ln g(x^i; w) \\ &= \sum_{i=1}^N y^i \ln \frac{1 - g(x^i; w)}{g(x^i; w)} + \ln g(x^i; w) \\ &= \sum_{i=1}^N y^i w^T x^i - \ln(1 + e^{w^T x^i}) \end{aligned}$$

Maximum likelihood estimation

$$\begin{aligned}\frac{\partial}{\partial w_j} l(w) &= \frac{\partial}{\partial w_j} \sum_{i=1}^N \{y^i w^T x^i - \ln(1 + e^{w^T x^i})\} \\ &= \sum_{i=1}^N x_j^i \{y^i - (1 - g(x^i; w))\} \\ &= \sum_{i=1}^N x_j^i \{y^i - p(y^i = 1 | x; w)\}\end{aligned}$$

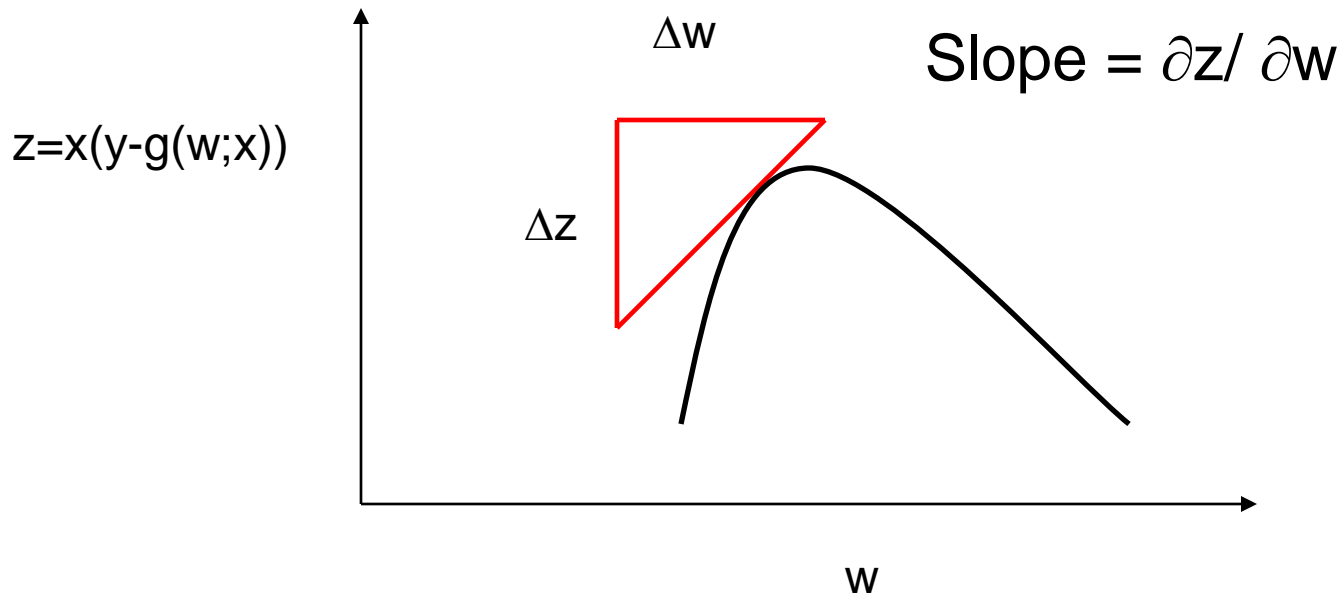
$$g(x; w) = \frac{1}{1 + e^{w^T x}}$$

$$1 - g(x; w) = \frac{e^{w^T x}}{1 + e^{w^T x}}$$

Bad news: No close form solution!

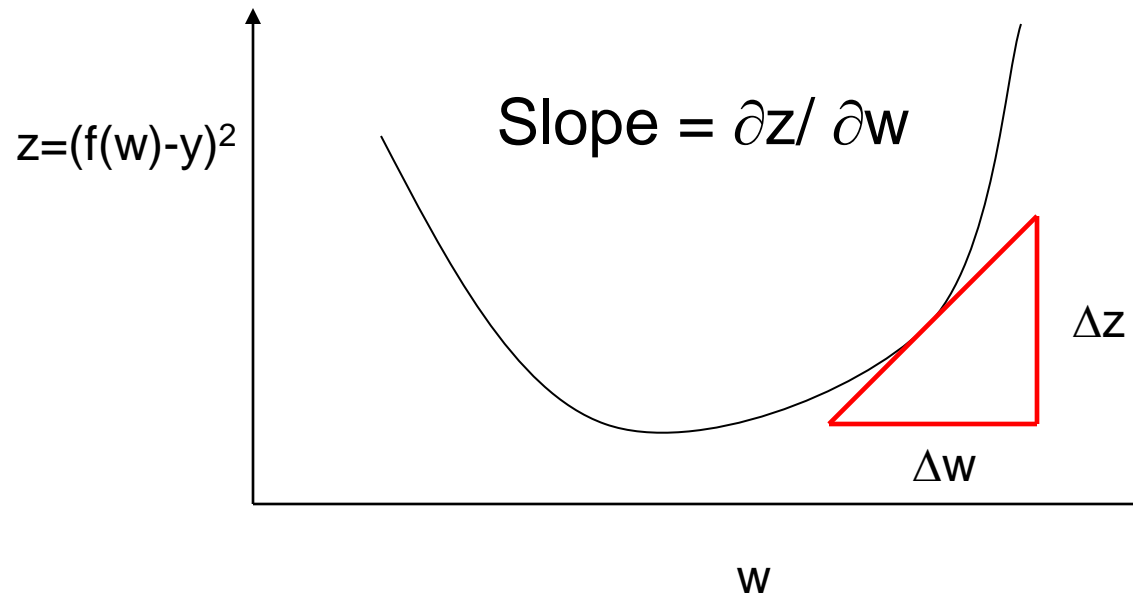
Good news: Concave function

Gradient ascent



- Going in the direction to the slope will lead to a larger z
- But not too much, otherwise we would go beyond the optimal w

Gradient descent



- Going in the *opposite* direction to the slope will lead to a smaller z
- But not too much, otherwise we would go beyond the optimal w

Gradient ascent for logistic regression

$$\frac{\partial}{\partial w_j} l(w) = \sum_{i=1}^N x_j^i \{y^i - (1 - g(x^i; w))\}$$

We use the gradient to adjust the value of w :

$$w_j \leftarrow w_j + \varepsilon \sum_{i=1}^N x_j^i \{y^i - (1 - g(x^i; w))\}$$

Where ε is a (small) constant

Example

Algorithm for logistic regression

1. Chose λ

2. Start with a guess for \mathbf{w}

3. For all j set
$$w_j \leftarrow w_j + \varepsilon \sum_{i=1}^N x_j^i \{y^i - (1 - g(x^i; \mathbf{w}))\}$$

4. If no improvement for $\sum_{i=1}^n (y^i - (1 - g(x^i; \mathbf{w})))^2$

stop. Otherwise go to step 3

Regularization

- Like with other data estimation problems, we may not have enough data to learn good models
- One way to overcome this is to ‘regularize’ the model, impose additional constraints on the parameters we are fitting.
- For example, lets assume that w_i comes from a Gaussian distribution with mean 0 and variance λ (where λ is a user defined parameter): $w_i \sim N(0, \lambda)$
- In that case we have:

$$p(y = 1, \theta | x) \propto p(y = 1 | x; \theta) p(\theta)$$

Regularization

- If we regularize the parameters we need to take the prior into account when computing the posterior for our parameters

$$p(y = 1, \theta | x) \propto p(y = 1 | x; \theta) p(\theta)$$

- Here we use a Gaussian model for the prior.
- Thus, the log likelihood changes to :

$$LL(y; w | x) = \sum_{i=1}^N y^i w^T x^i - \ln(1 + e^{w^T x^i}) - \sum_j \frac{w_j^2}{2\lambda}$$

After removing terms that are not dependent on w

- And the new update rule (after taking the derivative w.r.t. w_j) is:

$$w_j \leftarrow w_j + \varepsilon \sum_{i=1}^N x_j^i \{y^i - (1 - g(x^i; w))\} - \varepsilon \frac{w_j}{\lambda}$$

Also known as the MAP estimate

The variance of our prior model

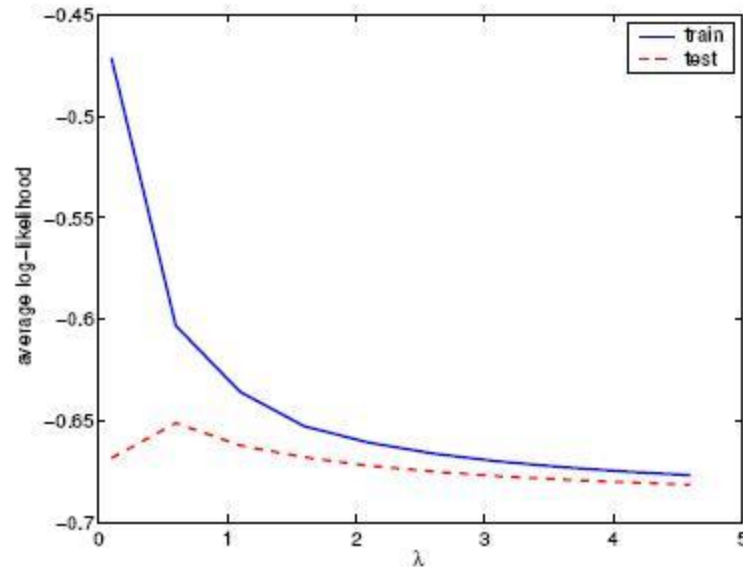
Regularization

- There are many other ways to regularize logistic regression
- The Gaussian model leads to an L2 regularization (we are trying to minimize the square of w)
- Another popular regularization is an L1 which tries to minimize $|w|$
- This often leads to many w_j 's being 0 resulting in compact models

The importance of the regularization parameter

- Too small does not have a big impact
- Too large overrides the data
- An example of the training/test conditional log likelihoods as a function of the regularization parameter λ

Average log
likelihood for data
only →



Logistic regression for more than 2 classes

- Logistic regression can be used to classify data from more than 2 classes:
- for $i < k$ we set

$$p(y = i | x; \theta) = g(w_{i0} + w_{i1}x_1 + \dots + w_{id}x_d) = g(\mathbf{w}_i^T \mathbf{x})$$

where

$$g(z_i) = \frac{e^{z_i}}{1 + \sum_{j=1}^{k-1} e^{z_j}} \quad z_i = w_{i0} + w_{i1}x_1 + \dots + w_{id}x_d$$

And for k we have

$$p(y = k | x; \theta) = 1 - \sum_{i=1}^{k-1} p(y = i | x; \theta) \Rightarrow$$
$$p(y = k | x; \theta) = \frac{1}{1 + \sum_{j=1}^{k-1} e^{z_j}}$$

Logistic regression for more than 2 classes

- Logistic regression can be used to classify data from more than 2 classes:
- for $i < k$ we set

$$p(y = i | x; \theta) = g(w_{i0} + w_{i1}x_1 + \dots + w_{id}x_d) = g(\mathbf{w}_i^T \mathbf{x})$$

where $g(z_i) = \frac{e^{z_i}}{1 + \sum_{j=1}^{k-1} e^{z_j}}$ $\leftarrow z_i = w_{i0} +$ Binary logistic regression is a special case of this rule

And for k we have $p(y = k | x; \theta) = 1 - \sum_{i=1}^{k-1} p(y = i | x; \theta) \Rightarrow$

$$p(y = k | x; \theta) = \frac{1}{1 + \sum_{j=1}^{k-1} e^{z_j}}$$

Update rule for logistic regression with multiple classes

$$\frac{\partial}{\partial w_{m,j}} l(w) = \sum_{i=1}^N x_j^i \{ \delta_m(y^i) - p(y^i = m | x^i; w) \}$$

Where $\delta(y^i)=1$ if $y^i=m$
and $\delta(y^i)=0$ otherwise

The update rule becomes:

$$w_{m,j} \leftarrow w_{m,j} + \varepsilon \sum_{i=1}^N x_j^i \{ \delta_m(y^i) - p(y^i = m | x^i; w) \}$$

Additive models

- Similar to what we did with linear regression we can extend logistic regression to other transformations of the data

$$p(y = 1 | x; w) = g(w_{i_0} + w_1\phi_1(x) + \dots + w_d\phi_d(x))$$

- As before, we are free to choose the basis functions

Important points

- Advantage of logistic regression over linear regression for classification
- Sigmoid function
- Gradient ascent / descent
- Regularization
- Logistic regression for multiple classes

Logistic regression

- The name comes from the **logit** transformation:

$$\log \frac{p(y = i | x; \theta)}{p(y = k | x; \theta)} = \log \frac{g(z_i)}{g(z_k)} = w_{i0} + w_{i1}x_1 + \dots + w_{id}x_d$$