

# Regression and Logistic Regression

# Regression

## ▶ Linear regression

▶ Objective : predict real values

▶ Training set

▶  $(\mathbf{x}^1, y^1), \dots, (\mathbf{x}^N, y^N)$

▶  $\mathbf{x} \in \mathbb{R}^n, y \in \mathbb{R}$  : single output regression

▶ Linear model

▶  $F(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} = \sum_{i=0}^n w_i x_i$  with  $x_0 = 1$

▶ Loss function

▶ Mean square error

□  $C = \frac{1}{2} \sum_{i=1}^N (y^i - \mathbf{w} \cdot \mathbf{x}^i)^2$

▶ Steepest descent gradient (batch)

▶  $\mathbf{w} = \mathbf{w}(t) - \epsilon \nabla_{\mathbf{w}} C, \nabla_{\mathbf{w}} C = \left( \frac{\partial C}{\partial w_1}, \dots, \frac{\partial C}{\partial w_n} \right)^T$

▶  $\frac{\partial C}{\partial w_k} = \frac{1}{2} \sum_{i=1}^N \frac{\partial}{\partial w_k} (y^i - \mathbf{w} \cdot \mathbf{x}^i)^2 = - \sum_{i=1}^N (y^i - \mathbf{w} \cdot \mathbf{x}^i) x_k^i$

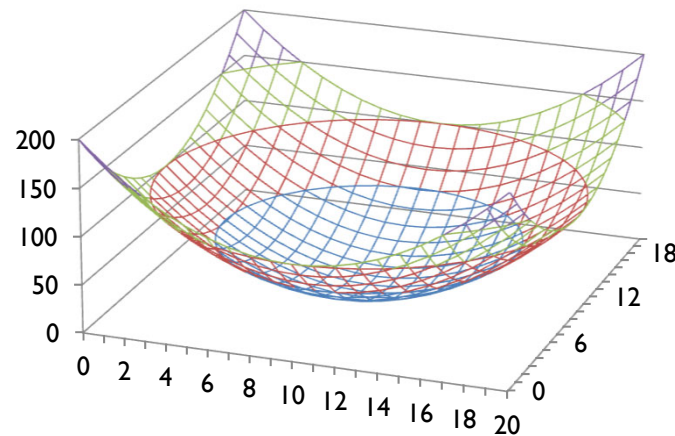
for component  $w_k$

▶  $\mathbf{w} = \mathbf{w}(t) + \epsilon \sum_{i=1}^N (y^i - \mathbf{w} \cdot \mathbf{x}^i) \mathbf{x}^i$

in vector form

## Regression

- ▶ Geometry of mean squares



- ▶ Regression with multiple outputs  $\mathbf{y} \in \mathbb{R}^p$ 
  - ▶ Simple extension:  $p$  independent linear regressions

## Probabilistic Interpretation

- ▶ Statistical model of linear regression
  - ▶  $y = \mathbf{w} \cdot \mathbf{x} + \epsilon$ , where  $\epsilon$  is a random variable (error term)
  - ▶ Hypothesis  $\epsilon$  is i.i.d. Gaussian
    - ▶  $\epsilon \sim N(0, \sigma^2)$ ,  $p(\epsilon) = \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{\epsilon^2}{2\sigma^2})$
    - ▶ The posterior distribution of  $y$  is then
    - ▶  $p(y | \mathbf{x}; \mathbf{w}) = \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(y - \mathbf{w} \cdot \mathbf{x})^2}{2\sigma^2})$
  - ▶ Likelihood
    - ▶  $L(\mathbf{w}) = \prod_{i=1}^N p(y^i | \mathbf{x}^i; \mathbf{w})$ 
      - Likelihood is a function of  $\mathbf{w}$ , it is computed on the training set
  - ▶ Maximum likelihood principle
    - ▶ Choose the parameters  $\mathbf{w}$  maximizing  $L(\mathbf{w})$  or any increasing function of  $L(\mathbf{w})$
  - ▶ In practice, one optimizes the log likelihood  $l(\mathbf{w}) = \log L(\mathbf{w})$ 
    - ▶  $l(\mathbf{w}) = N \log \left( \frac{1}{\sqrt{2\pi}\sigma} \right) - \frac{1}{2\sigma^2} \sum_{i=1}^N (y^i - \mathbf{w} \cdot \mathbf{x}^i)^2$
    - ▶ This is the MSE criterion
- ▶ This provides a probabilistic interpretation of regression
  - ▶ Under a gaussian hypothesis max likelihood is equivalent to MSE minimization

## Logistic regression – 2 classes

- ▶ Linear regression can be used (in practice) for regression or classification
- ▶ For classification a proper model is logistic regression

- ▶  $F_w(x) = \sigma(w \cdot x) = \frac{1}{1 + \exp(-w \cdot x)}$

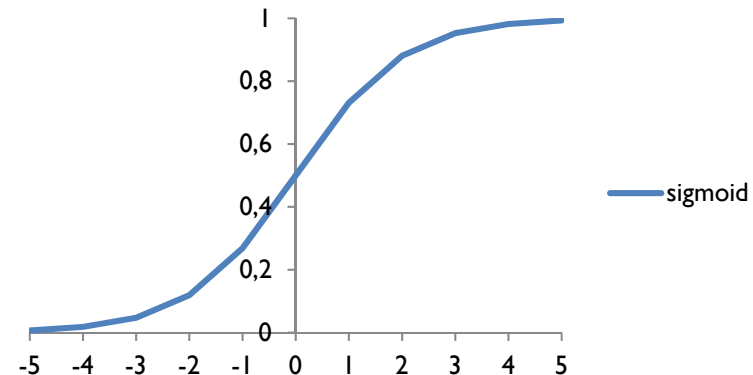
- ▶ Logistic (or sigmoid) function

- ▶  $\sigma(z) = \frac{1}{1 + \exp(-z)}$

- hint

- $\sigma'(z) = \sigma(z)(1 - \sigma(z))$

- ▶ Hyp:  $y \in \{0,1\}$



## Logistic regression – 2 classes

### Probabilistic interpretation

- ▶ Since  $y \in \{0,1\}$ , we make a Bernoulli hypothesis for the posterior distribution
  - ▶  $p(y = 1|\mathbf{x}; \mathbf{w}) = F_w(\mathbf{x})$  et  $p(y = 0|\mathbf{x}; \mathbf{w}) = 1 - F_w(\mathbf{x})$
  - ▶ In compact format
    - $p(y|\mathbf{x}; \mathbf{w}) = (F_w(\mathbf{x}))^y (1 - F_w(\mathbf{x}))^{1-y}$  with  $y \in \{0,1\}$
- ▶ Likelihood
  - ▶  $L(\mathbf{w}) = \prod_{i=1}^N (F_w(\mathbf{x}^i))^{y^i} (1 - F_w(\mathbf{x}^i))^{1-y^i}$
- ▶ Log-likelihood
  - ▶  $l(\mathbf{w}) = \sum_{i=1}^N y^i \log F_w(\mathbf{x}^i) + (1 - y^i) (\log(1 - F_w(\mathbf{x}^i)))$ 
    - This is minus the cross-entropy between the target and the estimated posterior distribution
  - ▶ Steepest descent algorithm (batch) for minimizing cross entropy
    - ▶ Componentwise:  $\frac{\partial l(\mathbf{w})}{\partial w_k} = \sum_{i=1}^N (y^i - F_w(\mathbf{x}^i)) x_k^i$
    - ▶ **Vector** form:  $\nabla_{\mathbf{w}} l = \sum_{i=1}^N (y^i - F_w(\mathbf{x}^i)) \mathbf{x}^i$
    - ▶ Algorithm
      - $\mathbf{w} = \mathbf{w} - \epsilon \nabla_{\mathbf{w}} \mathcal{C} = \mathbf{w} + \epsilon \sum_{i=1}^N (y^i - F_w(\mathbf{x}^i)) \mathbf{x}^i$

## Multivariate logistic regression

- ▶ Consider a  $p$  class classification problem
- ▶ Classes are encoded by “one hot” indicator vectors. Each vector is of dimension  $p$ 
  - ▶ Class 1:  $\mathbf{y} = (1, 0, \dots, 0)^T$
  - ▶ Class 2 :  $\mathbf{y} = (0, 1, \dots, 0)^T$
  - ▶ ...
  - ▶ Class  $p$ :  $\mathbf{y} = (0, 0, \dots, 1)^T$
- ▶  $F_{\mathbf{W}}(\mathbf{x})$  is a vector valued function with values in  $R^p$ 
  - ▶ Its component  $i$  is a **softmax function** (generalizes the sigmoid)
    - ▶  $\hat{y}_i = F_{\mathbf{W}}(\mathbf{x})_i = \frac{\exp(\mathbf{w}_i \cdot \mathbf{x})}{\sum_{j=1}^p \exp(\mathbf{w}_j \cdot \mathbf{x})}$ 
      - Note : here  $\mathbf{w}_j \in R^n$  is a vector,  $\hat{y}_i \in R$  is the  $i^{th}$  component of  $\hat{\mathbf{y}}$
- ▶ The probabilistic model for the posterior is a multinomial distribution
  - ▶  $p(\text{Class} = i | \mathbf{x}; \mathbf{w}) = \frac{\exp(\mathbf{w}_i \cdot \mathbf{x})}{\sum_{j=1}^p \exp(\mathbf{w}_j \cdot \mathbf{x})} = \text{softmax}(\mathbf{w}_i \cdot \mathbf{x})$

# Multivariate logistic regression

## ▶ Notations

- ▶  $\mathbf{s}^i = W\mathbf{x}^i$  is the logit for input  $\mathbf{x}^i$ 
  - ▶  $W = (\mathbf{w}_1, \dots, \mathbf{w}_p)^T$  is a  $p \times n$  matrix of weights
  - ▶  $\mathbf{s}^i = (s_1^i, \dots, s_p^i)^T \in \mathbb{R}^p$
- ▶  $\hat{\mathbf{y}}^i = \text{softmax}(\mathbf{s}^i)$  is the output for input  $\mathbf{x}^i$  (here  $\sigma$  applies component-wise, i.e.  $\hat{y}_j^i = \text{softmax}(s_j^i)$ )
  - ▶  $\hat{\mathbf{y}}^i = (\hat{y}_1^i, \dots, \hat{y}_p^i)^T \in \mathbb{R}^p$

## ▶ Let $\hat{\mathbf{y}}$ be a computed output for input $\mathbf{x}$ (we drop the index $i$ for simplicity), then

- ▶  $\frac{\partial \hat{y}_j}{\partial s_i} = \hat{y}_j(I_{ji} - \hat{y}_i)$  with  $I_{ji}$  elements of the identity matrix (1)

## ▶ Likelihood

- ▶  $L(W) = p(Y|X; W) = \prod_{i=1}^N \prod_{j=1}^p (\hat{y}_j^i)^{y_j^i}$ ,  $X$  and  $Y$  are the column wise matrices of input and output vector

## ▶ Log likelihood

- ▶  $l(W) = \sum_{i=1}^N \sum_{j=1}^p y_j^i \ln \hat{y}_j^i$  again this is minus the cross entropy for the multiclass classification problem

## ▶ Gradient of the log likelihood

- ▶  $\nabla_{w_k} l(W) = - \sum_{i=1}^N (\hat{y}_k^i - y_k^i) \mathbf{x}^i$  by using identity (1)

## ▶ Training algorithm

- ▶ As before, one may use a gradient method for maximizing the log likelihood.
- ▶ When the number of classes is large, computing the soft max is prohibitive, alternatives are required

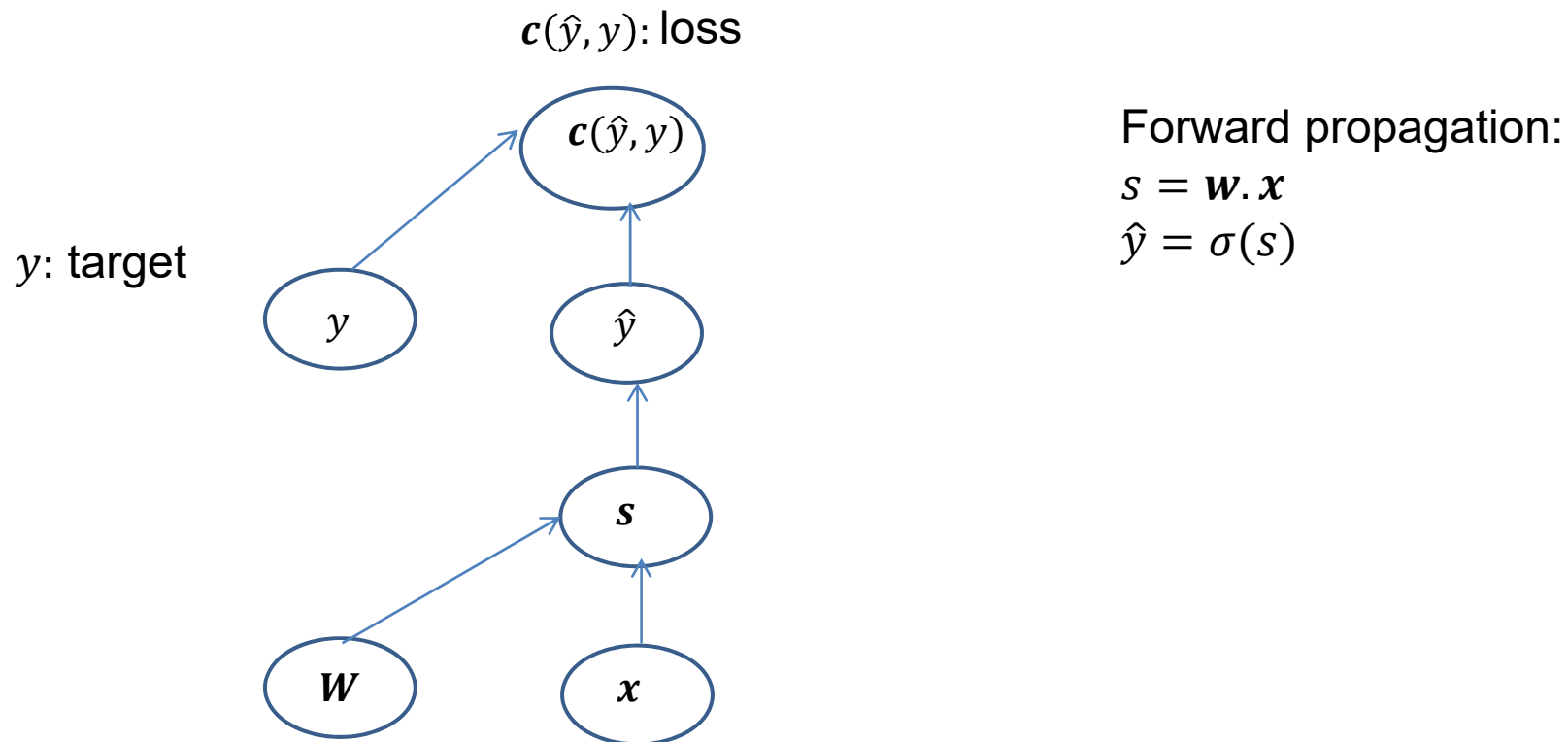


## Probabilistic interpretation for non linear models

- ▶ These results extend to non linear models, e.g. when  $F_w(x)$  is a NN
- ▶ Non linear regression
  - ▶ Max likelihood is equivalent to MSE loss optimization under the Gaussian hypothesis
    - ▶ For multivariate ( $y \in R, x \in R^n$ ) non linear regression we have
    - ▶  $y = F_w(x) + \epsilon, \epsilon \sim N(0, \sigma^2)$
    - ▶  $p(y | x; w) = \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(y-F(x))^2}{2\sigma^2})$
  - ▶ log – likelihood  $l(w)$ 
    - ▶  $l(w) = N \log \left( \frac{1}{\sqrt{2\pi}\sigma} \right) - \frac{1}{2\sigma^2} \sum_{i=1}^N (y^i - F(x^i))^2$
- ▶ Classification
  - ▶ Max likelihood is equivalent to cross entropy maximization under Bernoulli/multinomial distribution
    - 2 classes: if  $y$  is binary and we make the hypothesis that it is conditionally Bernoulli with probability  $F(x) = p(y = 1|x)$  we get the cross entropy loss
    - More than 2 classes: same as logistic regression with multiple outputs

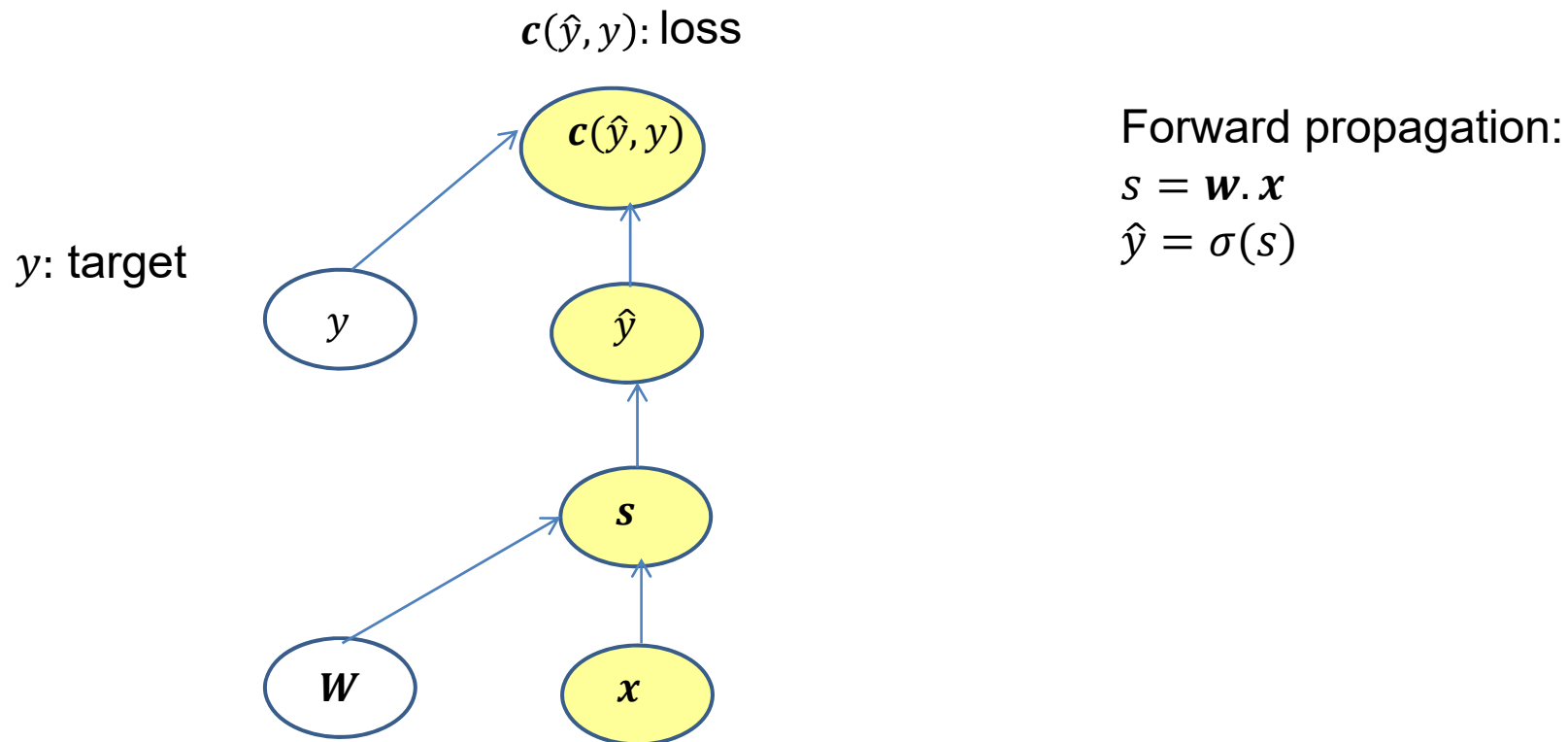
## Logistic regression – Computational graph -SGD

### ► Forward pass



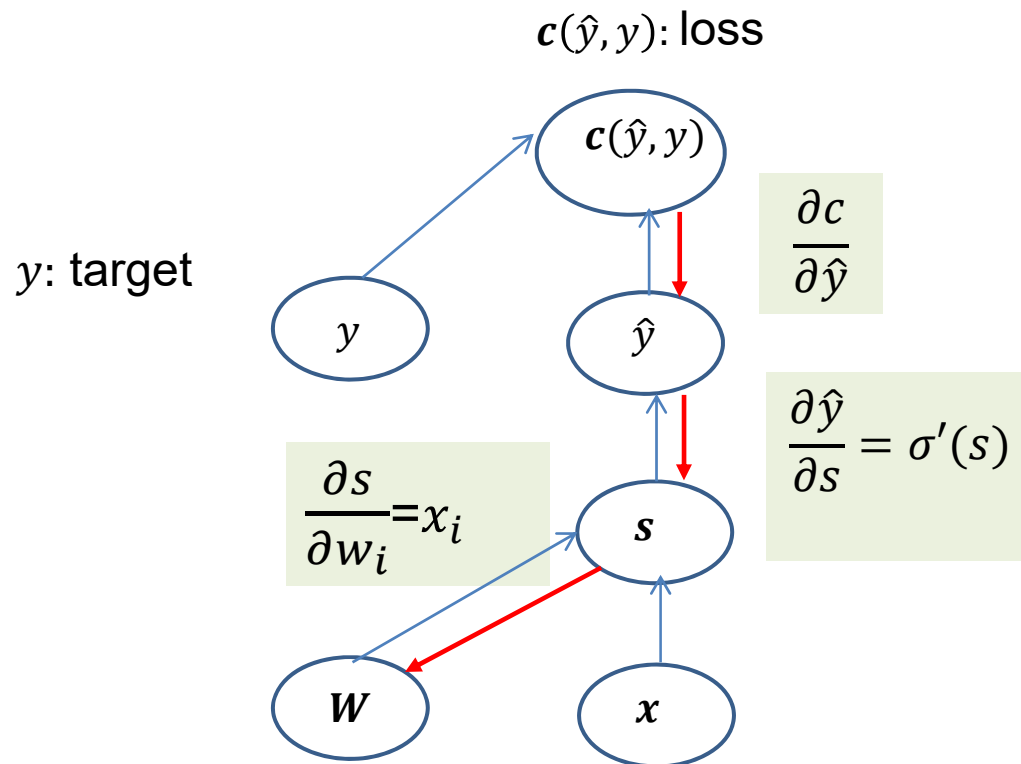
## Logistic regression – Computational graph - SGD

### ► Forward pass



## Logistic regression – Computational graph - SGD

### ► Backward pass



Backward propagation:

$$\frac{\partial c}{\partial s} = \frac{\partial c}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial s}$$
$$\frac{\partial c}{\partial w_i} = \frac{\partial c}{\partial s} \frac{\partial s}{\partial w_i}$$

$$\frac{\partial c}{\partial w_i} = \frac{\partial c}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial s} \frac{\partial s}{\partial w_i}$$

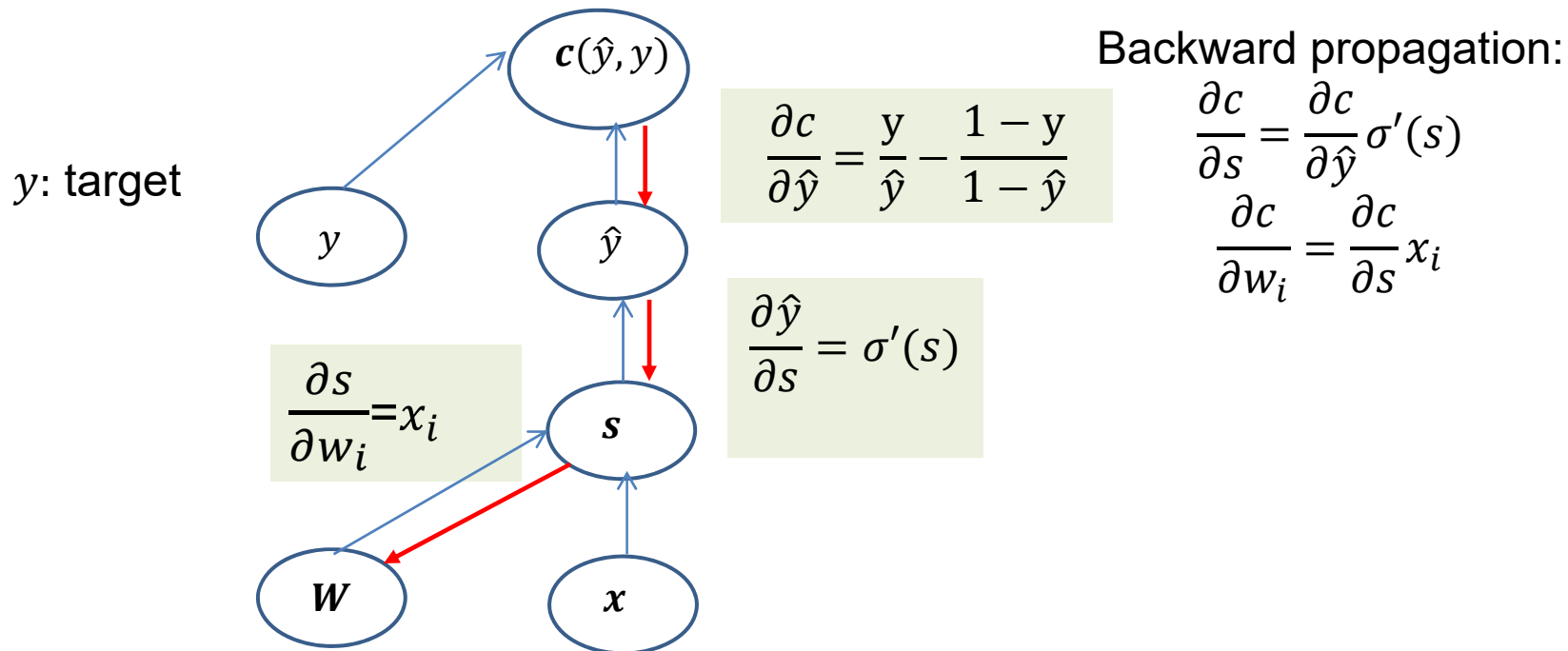
Chain Rule

## Logistic regression – Computational graph - SGD

- ▶ **Backward pass** For the cross entropy loss  

$$l(\mathbf{w}) = \sum_{i=1}^N y^i \log \hat{y}^i + (1 - y^i) \log(1 - \hat{y}^i) = \sum_{i=1}^N c(\hat{y}^i, y^i)$$

$c(\hat{y}, y)$ : loss



$$\frac{\partial c}{\partial w_i} = \left( \frac{y}{\hat{y}} - \frac{1-y}{1-\hat{y}} \right) \sigma'(s) x_i$$

# Probabilistic interpretation of NN outputs

## Mean Square loss

- ▶ Derived here for multivariate regression (1 output), trivial extension to multiple outputs
- ▶ Holds for any continuous functional (regression, logistic regression, NNs, etc)
- ▶ Risk  $R = E_{x,y} [(y - h(x))^2]$
- ▶ The minimum of  $R$ ,  $\text{Min}_h R$ , is obtained for  $h^*(x) = E_y[y|x]$
- ▶ The risk  $R$  pour the family of functions  $F_w(x)$  decomposes as follows:
  - ▶  $R = E_{x,y} [(y - F_w(x))^2]$
  - ▶  $R = E_{x,y} [(y - E_y[y|x])^2] + E_{x,y} [(E_y[y|x] - F_w(x))^2]$
- ▶ Let us consider  $E_y [(y - E_y[y|x])^2]$ 
  - ▶ This term is independent of the model  $F_w(\cdot)$  and only depends on the problem characteristics (the data distribution).
  - ▶ It represents the min error that could be obtained for this data distribution
  - ▶  $h^*(x) = E_y[y|x]$  is the optimal solution to  $\text{Min}_h R$
- ▶ Minimizing  $E_{x,y} [(y - F_w(x))^2]$  is equivalent to minimizing  $E_{x,y} [(E_y[y|x] - F_w(x))^2]$ 
  - ▶ The optimal solution  $F_{w*}(x) = \text{argmin}_w E_{x,y} [(E_y[y|x] - F_w(x))^2]$  is the best mean square approximation of  $E[y|x]$

## Probabilistic interpretation of NN outputs

### ► Classification

- Let us consider multi-class classification with one hot encoding of the target outputs
  - i.e.  $\mathbf{y} = (0, \dots, 0, 1, 0, \dots, 0)^T$  with a 1 at position  $i$  if the target is class  $i$  and zero everywhere else
  - $h_i^* = E_y[y|x] = 1 * P(C_i|x) + 0 * (1 - P(C_i|x)) = P(C_i|x)$
  - i.e.  $F_{w^*}()$  is the best LMS approximation of the Bayes discriminant function (which is the optimal solution for classification with 0/1 loss)
- More generally with binary targets
  - $h_i^* = P(y_i = 1|x)$

### ► Note

- Similar results hold for the cross entropy criterion
- Precision on the computed outputs depends on the task
  - Classification: precision might not be so important (max decision rule, one wants the correct class to be ranked above all others)
  - Posterior probability estimation: precision is important