$$4/13/18$$

## Non-linear regression
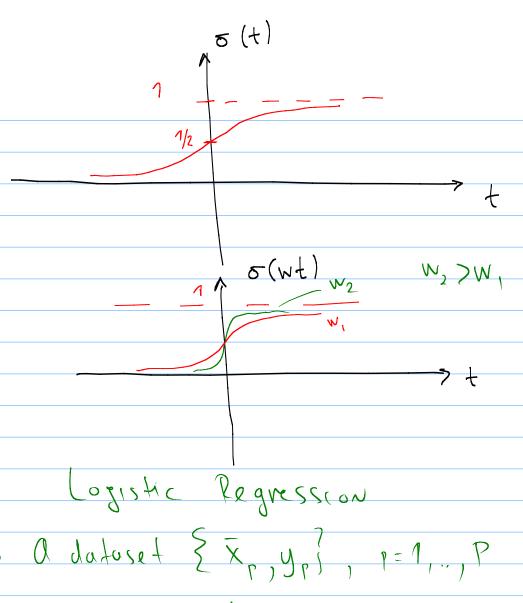
- Associated cost is nonlinear in its parameters

- Logistic regression

### Logistic sigmoid function

18th century   Verhulst

$f$ population          max capacity

growth rate $\dfrac{df}{dt} = f(1-f)$                     $f(0) = 1/2$

current population        remaining capacity

$$\sigma(t) = \dfrac{1}{1+e^{-t}}$$

$$\sigma'(t) = \dfrac{-(-e^{-t})}{(1+e^{-t})^2} = \underbrace{\dfrac{1}{1+e^{-t}}}_{\sigma(t)} \cdot \underbrace{\dfrac{e^{-t}}{1+e^{-t}}}_{1-\sigma(t)}$$

$$1 - \dfrac{1}{1+e^{-t}} = \dfrac{1+e^{-t}-1}{1+e^{-t}} = \dfrac{e^{-t}}{1+e^{-t}}$$

$$\sigma'(t) = \sigma(t)(1-\sigma(t))$$

$\sigma(t)$

$1$

$1/2$

$t$

$\sigma(wt)$  $w_2$   $w_2 > w_1$

$1$

$w_1$

$t$

# Logistic Regression

- A dataset $\{\bar{x}_p, y_p\}$, $p = 1, \ldots, P$

  is distributed sigmoidally

  $$\sigma(b + \bar{x}_p^T \bar{w}) \approx y_p, \quad p = 1, \ldots, P$$

- Least Squares approach

  $$g(b, \bar{w}) = \sum_{p=1}^{P} \left( \sigma(b + \bar{x}_p^T \bar{w}) - y_p \right)^2$$

- Compact Notation

  $$\tilde{x}_p = \begin{bmatrix} 1 \\ \bar{x}_p \end{bmatrix}_{(N+1) \times 1} \qquad \tilde{w} = \begin{bmatrix} b \\ \bar{w} \end{bmatrix}_{(N+1) \times 1}$$

  $$g(\tilde{w}) = \sum_{p=1}^{P} \left( \sigma(\tilde{x}_p^T \tilde{w}) - y_p \right)^2$$

$$\nabla_{\tilde{w}} g(\tilde{w}) = \sum_{p=1}^{P} \nabla \left( \sigma(\tilde{x}_p^T \tilde{w}) - y_p \right)^2$$

$$= \sum_{p=1}^{P} 2 \left( \sigma(\tilde{x}_p^T \tilde{w}) - y_p \right) \cdot \nabla_{\tilde{w}} \left( \sigma(\tilde{x}_p^T \tilde{w}) - y_p \right)$$

$$= \sum_p 2 \left( \sigma(\tilde{x}_p^T \tilde{w}) - y_p \right) \cdot \sigma(\tilde{x}_p^T \tilde{w})\left(1 - \sigma(\tilde{x}_p^T \tilde{w})\right)$$

$$\cdot \underbrace{\nabla(\tilde{x}_p^T \tilde{w})}_{\tilde{x}_p}$$

Gradient descent

$$\tilde{w}_{k+1} = \tilde{w}_k - \alpha_k \cdot \nabla g(\tilde{w}_k)$$

---

## Regularization

original cost $\quad \min_{\tilde{w}} g(\tilde{w})$

regularized cost $\quad \min_{\tilde{w}} \left( g(\tilde{w}) + \lambda R(\tilde{w}) \right)$ ← regularization parameter

Examples of $R(\tilde{w})$

- $R(\tilde{w}) = \| C\tilde{w} \|_2^2 \qquad$ incorporation of prior knowledge

- $R(\tilde{w}) = \| \tilde{w} \|_1^2 \qquad \ell_1$ norm enforces sparcity

- avoid overfitting $\quad R(\tilde{w}) = \| \tilde{w} \|_2^2$ ←

- convexify $g(\tilde{w}) \quad R(\tilde{w}) = \| \tilde{w} \|_2^2$