# Logistic Regression

Mohammad Emtiyaz Khan
EPFL

Oct 8, 2015
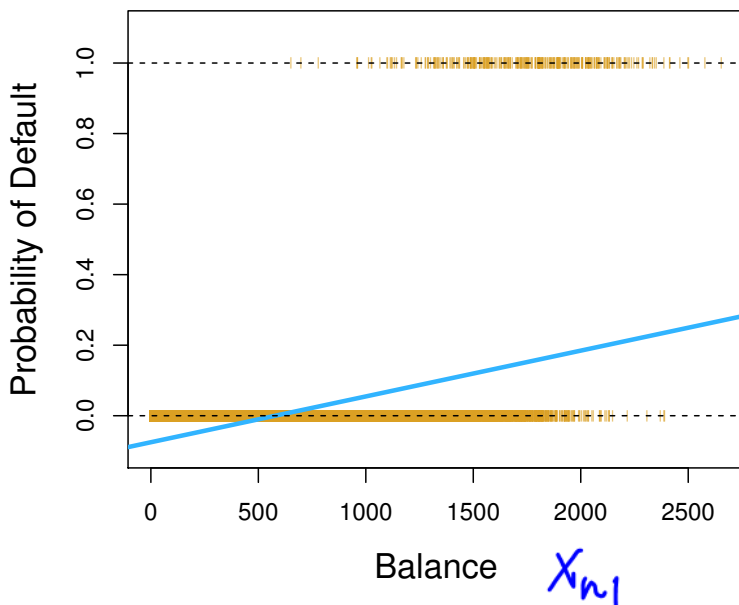


ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

# Classification with linear regression

We can use $y = 0$ for $\mathcal{C}_1$ and $y = 1$ for $\mathcal{C}_2$ (or vice-versa), and simply use least-squares to predict $\hat{y}_*$ given $\mathbf{x}_*$. We can predict $\mathcal{C}_1$ when $\hat{y}_* < 0.5$ and $\mathcal{C}_2$ when $\hat{y}_* > 0.5$.
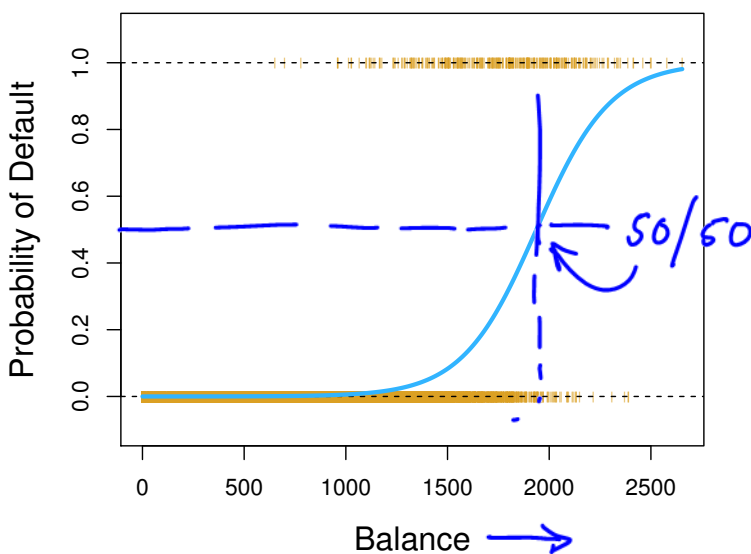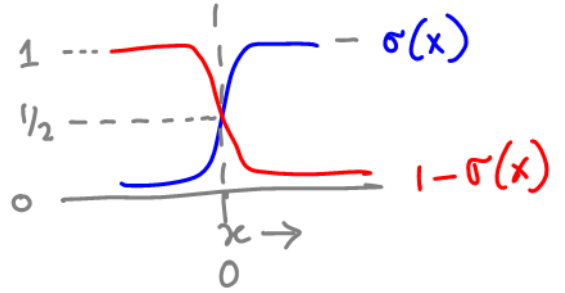
Default



$$p\left(\text{Default} = 1 \middle/ X\right)$$

$$\beta_0 + \beta_1 x_{n1} \simeq y_n$$

Any problems with this approach?



50/50

# Logistic regression

We need to model $p(y = C_1 | \mathbf{x})$ and $p(y = C_2 | \mathbf{x})$ such that they both are $> 0$ and also sum to 1. For a new input $\mathbf{x}_*$, we can classify to $C_1$ when $p(\hat{y}_* | \mathbf{x}_*) < 0.5$.

*(handwritten annotations: $\sigma(x)$, $1 - \sigma(x)$, $\sigma(x) + (1 - \sigma(x)) = 1$)*

We will use the logistic function.
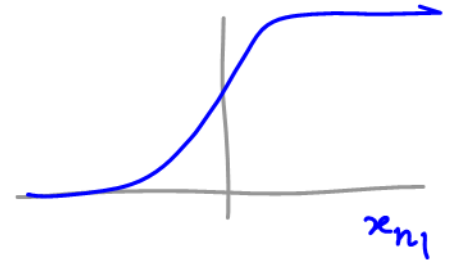
$$\sigma(x) = \frac{\exp(x)}{1 + \exp(x)}, \qquad 1 - \sigma(x) = \frac{1}{1 + \exp(x)}$$

We pass the linear-regression model $\eta_n = \widetilde{\mathbf{x}}^T \boldsymbol{\beta}$ through the logistic function to get the probabilities.

$$p(y_n = C_1 | \mathbf{x}_n) = \sigma(\eta_n), \qquad p(y_n = C_2 | \mathbf{x}_n) = 1 - \sigma(\eta_n)$$

This figure visualizes the probabilities obtained for a 2-D problem (taken from KPM Chapter 7).

*(handwritten annotations: $x_{n1}$, $\boldsymbol{\beta}^T \widetilde{\mathbf{x}}_n$)*

*(handwritten: ·Visualize)*



*(handwritten: Replace $\underline{w}$ by $\underline{\beta}$)*

2

# The probabilistic model

Assuming that each $y_n$ is independent of others, we can define the probability of **y** given **X** and $\boldsymbol{\beta}$:

(A) $p(\mathbf{y}|\mathbf{X}, \boldsymbol{\beta}) = \prod\limits_{n=1}^{N} p(y_n|\mathbf{x}_n)$ ← $\sigma\left(\boldsymbol{\beta}^T \widetilde{\underline{x}}_n\right)$ or $\left[1 - \sigma\left(\boldsymbol{\beta}^T \widetilde{\underline{x}}_n\right)\right]$

(B) $= \prod\limits_{n:y_n=\mathcal{C}_1} p(y_n = \overset{1}{\mathcal{C}_1}|\mathbf{x}_n) \prod\limits_{n:y_n=\mathcal{C}_2} p(y_n = \overset{0}{\mathcal{C}_2}|\mathbf{x}_n)$

$\underbrace{\phantom{n:y_n=\mathcal{C}_1}}\ \sigma\left(\underline{x}_n^T \boldsymbol{\beta}\right) \quad \underbrace{\phantom{n:y_n=\mathcal{C}_2}}\ \left(1 - \sigma\left(\widetilde{\underline{x}}_n^T \boldsymbol{\beta}\right)\right)$

A better way to write this is to use the coding $y_n \in \{0, 1\}$.

$\begin{cases} \phi\left(y_n = \mathcal{C}_1 / \underline{x}_n\right) \overset{\text{when}}{\underset{\cup}{}} y_n = \mathcal{C}_1 \\ \phi\left(y_n = \mathcal{C}_2 / \underline{x}_n\right) \overset{\text{when}}{\underset{0}{}} y_n = \underline{\mathcal{C}_2} \end{cases}$

(C) $p(\mathbf{y}|\mathbf{X}, \boldsymbol{\beta}) = \prod\limits_{n=1}^{N} \sigma(\eta_n)^{y_n}[1 - \sigma(\eta_n)]^{1-y_n}$ ‑ ‑ ‑ ‑ ‑ ‑

The <u>log-likelihood</u> is given as follows:

$\log \phi\left(\underline{y}/X, \boldsymbol{\beta}\right)$

(D) $\mathcal{L}_{mle}(\boldsymbol{\beta}) = \sum\limits_{n=1}^{N} y_n \log \sigma(\widetilde{\mathbf{x}}_n^T \boldsymbol{\beta}) + (1 - y_n) \log[1 - \sigma(\widetilde{\mathbf{x}}_n^T \boldsymbol{\beta})]$

(E) $= \sum\limits_{n=1}^{N} y_n \widetilde{\mathbf{x}}_n^T \boldsymbol{\beta} - \log[1 + \exp(\widetilde{\mathbf{x}}_n^T \boldsymbol{\beta})]$

↑ Gives us (D)

Proof: (C) → (E)

$\log \sigma\left(\widetilde{x}_n^T \boldsymbol{\beta}\right) = \log \dfrac{e^{\widetilde{x}_n^T \boldsymbol{\beta}}}{1 + e^{\widetilde{x}_n^T \boldsymbol{\beta}}} = \widetilde{x}_n^T \boldsymbol{\beta} - \log\left[1 + e^{\widetilde{x}_n^T \boldsymbol{\beta}}\right]$

$\log\left[1 - \sigma()\right] = \log \dfrac{1}{1 + e^{[\ ]}} = -\log\left[1 + e^{\widetilde{x}_n^T \boldsymbol{\beta}}\right]$

To get (E), Add these for all $b$     $\left(y_n + 1 - y_n\right) = 1$ ↘

$= \sum\limits_{n=1}^{N} y_n \widetilde{x}_n^T \boldsymbol{\beta} - y_n \log\left[1 + \exp(\widetilde{x}_n^T \boldsymbol{\beta})\right] - (1 - y_n)\log\left[\quad\right]$

3

$$\mathcal{L}(\beta) := \sum_n y_n \widetilde{x}_n^T \beta - \log\left[1 + \exp(\widetilde{x}_n^T \beta)\right]$$

# Maximum likelihood

We will use the following fact to derive the gradient.

$$\frac{\partial}{\partial t} \log\left[1 + e^t\right] = \frac{e^t}{1 + e^t}$$
$$= \sigma(t)$$

$$\frac{\partial}{\partial x} \log[1 + \exp(x)] = \sigma(x)$$

Taking the gradient of the log-likelihood, we get the following:

$$\frac{\partial \mathcal{L}}{\partial \beta} = \sum_n \widetilde{x}_n y_n - \widetilde{x}_n \sigma(\widetilde{x}_n^T \beta)$$

$$= \sum_n \widetilde{x}_n \left[y_n - \sigma(\widetilde{x}_n^T \beta)\right]$$

$$\mathbf{g} := \frac{\partial \mathcal{L}}{\partial \boldsymbol{\beta}} = \widetilde{\mathbf{X}}^T [\sigma(\widetilde{\mathbf{X}}\boldsymbol{\beta}) - \mathbf{y}] \quad = 0$$

Least-squares: $-\widetilde{X}^T[\widetilde{X}\beta - y]$

This is similar to the normal equation for least-squares.

Define $\sigma(\widetilde{X}\beta) := \begin{bmatrix} \sigma(\widetilde{x}_1^T \beta) \\ \sigma(\widetilde{x}_2^T \beta) \\ \vdots \\ \sigma(\widetilde{x}_N^T \beta) \end{bmatrix}$

There are no closed-form solutions, but we can use gradient descent.

# Convexity

The negative of the log-likelihood $-\mathcal{L}_{mle}(\boldsymbol{\beta})$ is convex.

$$\max_\beta \mathcal{L}(\beta)$$
$$= \min_\beta -\mathcal{L}(\beta)$$

Proof I: The sum of a linear function and a (strictly) convex function is (strictly) convex.

$$\mathcal{L}(\beta) = \sum_{n=1}^{N} (\widetilde{x}_n^T \beta) y_n - \log(1 + e^{\widetilde{x}_n^T \beta})$$

$$-\mathcal{L}(\beta) = \sum_{n=1}^{N} -(\widetilde{x}_n^T \beta) + \log(1 + e^{\widetilde{x}_n^T \beta})$$

Proof II: The Hessian of a convex function is positive semi-definite and for a strictly-convex function it is positive definite.

$$-\frac{\partial \mathcal{L}}{\partial \beta} = \sum_n -y_n \widetilde{x}_n + \sigma(\widetilde{x}_n^T \beta) \widetilde{x}_n$$

$$-\frac{\partial^2 \mathcal{L}}{\partial \beta \partial \beta^T} = \frac{\partial}{\partial \beta \partial \beta^T} \sigma(\widetilde{x}_n^T \beta) \widetilde{x}_n$$

Q

4

# Hessian of the Log-Likelihood

We will use the following fact:

$$\frac{\partial \sigma(t)}{\partial t} = \sigma(t)[1 - \sigma(t)]$$

$$\frac{\partial}{\partial t} \frac{e^t}{1+e^t} = -\frac{e^{2t}}{(1+e^t)^2} + \frac{e^t}{1+e^t}$$

$$= -\left[\sigma(t)\right]^2 + \sigma(t)$$

$$= \sigma(t)\left[1 - \sigma(t)\right]$$

Taking the derivative of the gradient we get the Hessian,

$$\mathbf{H}(\boldsymbol{\beta}) := -\frac{\partial \mathbf{g}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}^T} = \widetilde{\mathbf{X}}^T \mathbf{S} \widetilde{\mathbf{X}}$$

$$\widetilde{X}^T \widetilde{X}$$

where $\mathbf{S}$ is a $N \times N$ diagonal matrix with diagonals

$$S_{nn} = \sigma(\widetilde{\mathbf{x}}_n^T \boldsymbol{\beta})[1 - \sigma(\widetilde{\mathbf{x}}_n^T \boldsymbol{\beta})]. \leftarrow$$

$$\frac{\partial^2 \mathcal{L}}{\partial \underline{\beta} \, \partial \underline{\beta}^T} = \sum_{n=1}^{N} \frac{\partial}{\partial \underline{\beta}^T} \sigma\left(\underline{\widetilde{x}}_n^T \underline{\beta}\right) \underbrace{\underline{\widetilde{x}}_n}_{(D+1)}$$

$$(D+1)\times(D+1)$$

$$= \sum_{n=1}^{N} \underbrace{\underline{\widetilde{x}}_n}_{(D+1)\times 1} \sigma\left(\underline{\widetilde{x}}_n^T \underline{\beta}\right) \underbrace{\left[1 - \sigma(\underline{\widetilde{x}}_n^T \underline{\beta})\right]}_{1\times 1} \underbrace{\underline{\widetilde{x}}_n^T}_{1\times(D+1)}$$

$$\underbrace{\qquad}_{S_{nn}}$$

$$S = \begin{bmatrix} S_{11} & & & \\ & S_{22} & & \\ & & \ddots & \\ & & & S_{NN} \end{bmatrix}_{N\times N}$$

Is the negative of the log-likelihood *strictly* convex?   No

$$\underline{a}^T \mathbf{H} \, \underline{a} > 0$$

$$\Rightarrow \underline{a}^T \widetilde{X}^T S \widetilde{X} \underline{a} > 0$$

$$\Rightarrow \underline{a} \widetilde{X}^T S^{V_2} S^{V_2} \widetilde{X} \underline{a} > 0$$

$$\underbrace{\qquad}_{t\rightarrow}$$

$$\Rightarrow \underline{t}^T \underline{t} > 0$$

but when
$\exists \, \underline{a}$ s.t $\widetilde{X}\underline{a} = 0$
$\neq 0$
then $\underline{t}^T \underline{t} = 0$

$$\frac{\partial^2 \mathcal{L}}{\partial \underline{\beta} \, \partial \underline{\beta}^T} = \widetilde{X}^T S \widetilde{X}$$

else

# Newton's Method

$$\underline{g}(\underline{\beta}) \stackrel{\Delta}{=} -\frac{\partial \mathcal{L}}{\partial \underline{\beta}} = \widetilde{X}^T\left[\sigma(\widetilde{X}\underline{\beta}) - \underline{y}\right]$$

Gradient descent uses only first-order information and takes steps in the direction of the gradient.

gradient descent: $\underline{\beta}^{(k+1)} = \underline{\beta}^{(k)} - \alpha_k \underline{g}_k$

Newton: $\underline{\beta}^{(k+1)} = \underline{\beta}^{(k)} - \alpha_k \mathbf{H}_k^{-1}\underline{g}_k$

Newton's method uses second-order information and takes steps in the direction that minimizes a quadratic approximation.

where $\underline{g}_k = \frac{\partial}{\partial \underline{\beta}}\left[-\mathcal{L}(\underline{\beta})\right]\Big|_{\underline{\beta}=\underline{\beta}^{(k)}}$

$\mathbf{H}_k = \frac{\partial^2}{\partial\underline{\beta}\,\partial\underline{\beta}^T}\Big|_{\underline{\beta}=\underline{\beta}^{(k)}}$

$$\boldsymbol{\beta}^{(k+1)} = \boldsymbol{\beta}^{(k)} - \alpha_k \mathbf{H}_k^{-1}\mathbf{g}_k$$

where $\mathbf{g}_k$ is the gradient.

Newton minimizes the quadratic approx. at $\beta^k$

$\mathcal{L}(\underline{\beta}) \cong \mathcal{L}(\underline{\beta}^{(k)}) + \underline{g}_k^T(\underline{\beta} - \underline{\beta}^{(k)})$

# Computational complexity

Compare the computational complexity of least-squares and Newton's method.

$+ (\underline{\beta} - \underline{\beta}^{(k)})^T \mathbf{H}_k (\underline{\beta} - \underline{\beta}^{(k)})$

$= Q(\underline{\beta}, \underline{\beta}^{(k)})$

$\underline{\beta}^* = \underset{\underline{\beta}}{\arg\min}\, Q(\underline{\beta}, \underline{\beta}^{(k)})$

$= \mathbf{H}_k^{-1}\underline{g}_k$

Newton's method is equivalent to solving many least-squares problems.

only applies to linear reg

| | |
|---|---|
| Gradient descent | $O(NDI)$ |
| least-squares | $O(ND^2 + D^3)$ |
| Newton method | $O[(ND^2 + D^3)I]$ |

$\left(\widetilde{X}^T \overset{(k)}{S} \widetilde{X}\right)^{-1} \widetilde{X}^T\left(\sigma(\widetilde{X}\underline{\beta}^{(k)}) - \underline{y}\right)$

$\underbrace{\qquad}_{\mathbf{H}_k} \qquad \underbrace{\qquad}_{\underline{g}_k}$

$D\times N \quad N\times D$

$D^2 N \qquad\qquad D^3$

Least-squares: $\left(\widetilde{X}^T\widetilde{X}\right)^{-1}\widetilde{X}^T(\widetilde{X}\underline{\beta} - \underline{y}) = 0$

# Penalized Logistic Regression

The cost-function can be unbounded when the data is linearly separable.

For a well-defined problem, we will regularize.

$$X^T S X + \lambda \begin{bmatrix} 0 & \\ & I_D \end{bmatrix}_D$$

$$\min_{\beta} - \sum_{n=1}^{N} \log p(y_n | \mathbf{x}_n^T \boldsymbol{\beta}) + \lambda \sum_{d=1}^{D} \beta_d^2$$

↑

$$\begin{bmatrix} \text{Regularization} \\ \text{Solves this} \\ \text{problem} \end{bmatrix}$$

In "linearly Separable" case,
as $\beta \to \infty$, $\mathcal{L}(\beta) \to \infty$.
Therefore the global minimum does not exist!

← $\beta \to \infty$

$$\begin{bmatrix} \text{Typically, a linearly sepa-} \\ \text{rable problem has many} \\ \text{Solutions.} \end{bmatrix}$$

Linearly separable in 1-D

0.6
0.5
0.4

$x \to$

$x < \tau$
⇒ $y = '+'$

$x > \tau \Rightarrow y = '0'$

"Margin" in 2-D

0.5

$x_2$

$x_1$

7

# Additional notes

## Derivation of Newton's method

The second-order approximation of a function is given as follows:

$$\mathcal{L}_Q(\boldsymbol{\beta}) := \mathcal{L}(\boldsymbol{\beta}^{(k)}) + \mathbf{g}_k^T(\boldsymbol{\beta} - \boldsymbol{\beta}^{(k)}) + \tfrac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\beta}^{(k)})^T \mathbf{H}_k(\boldsymbol{\beta} - \boldsymbol{\beta}^{(k)})$$

The minimum of $\mathcal{L}_Q$ is at $\boldsymbol{\beta}^{(k)} - \mathbf{H}_k^T \mathbf{g}_k$. A conservative option is to take a small step in this direction using step-size $\alpha_k$, which is the step used in Newton's method.

Set $\alpha_k$ using line search, e.g. the Armijo rule. See Section 8.3.2 of Kevin Murphy's book. A good implementation can be found on page 29 of Bertsekas book "Non-linear programming".

## Iterative Recursive Least-Squares (IRLS)

(IRLS) expresses Newton's method with $\alpha_k = 1$ as a sequence of least-squares problems. Below is the derivation and pseudo code.

$$\boldsymbol{\beta}^{(k+1)} = \boldsymbol{\beta}^{(k)} - \alpha_k \mathbf{H}_k^{-1} \mathbf{g}_k \tag{1}$$

$$= \boldsymbol{\beta}^{(k)} - (\widetilde{\mathbf{X}}^T \mathbf{S}_k \widetilde{\mathbf{X}})^{-1} \widetilde{\mathbf{X}}^T(\boldsymbol{\sigma}_k - \mathbf{y}) \tag{2}$$

$$= (\widetilde{\mathbf{X}}^T \mathbf{S}_k \widetilde{\mathbf{X}})^{-1}[(\widetilde{\mathbf{X}}^T \mathbf{S}_k \widetilde{\mathbf{X}})\boldsymbol{\beta}^{(k)} - \widetilde{\mathbf{X}}^T(\boldsymbol{\sigma}_k - \mathbf{y})]$$

$$= (\widetilde{\mathbf{X}}^T \mathbf{S}_k \widetilde{\mathbf{X}})^{-1}\widetilde{\mathbf{X}}^T \mathbf{X}_k[\widetilde{\mathbf{X}}\boldsymbol{\beta}^{(k)} + \mathbf{S}_k^{-1}(\mathbf{y} - \boldsymbol{\sigma}_k)]$$

$$= (\widetilde{\mathbf{X}}^T \mathbf{S}_k \widetilde{\mathbf{X}})^{-1}\widetilde{\mathbf{X}}^T \mathbf{S}_k \mathbf{z}_k \tag{3}$$

where $\mathbf{z}_k = \widetilde{\mathbf{X}}\boldsymbol{\beta}^{(k)} + \mathbf{S}_k^{-1}(\mathbf{y} - \boldsymbol{\sigma}_k)$.

```
1 for k = 1:maxIters
2     sig = sigmoid(tX*beta);
3     s = sig.*(1-sig);
4     z = tX*beta + (y-sig)./s;
5     beta = weightedLeastSquares(z,tX,s);
6 end
```

# Quasi-Newton

Read about L-BFGS in Section 8.3.5 of Kevin Murphy's book. The key idea is to approximate **H** usign a diagonal and a low-rank matrix.

## To do

1. Practice to derive the cost function using maximum likelihood estimation.

2. Understand the normal equation.

3. Understand the interpretation of log-odds (JWHT Chapter 3).

4. Learn to prove convexity using the positive-definite property of the Hessian.

5. Implement Newton's method (part of next week's lab).

6. Understand the relationship of Newton's Method with IRLS.

7. Do exercise 8.3 to 8.7 in KPM book. (Chapter on logistic Regression)