

Big Data Analytics

The Linear Model

Souhaib Ben Taieb

March 25, 2021

University of Mons

Table of contents

Linear Classification

Linear Regression

Important notions in nonlinear optimization

Back to linear regression

Logistic Regression

Table of contents

Linear Classification

Linear Regression

Important notions in nonlinear optimization

Back to linear regression

Logistic Regression

A real data set



Input representation

'raw' input $\mathbf{x} = (x_0, x_1, x_2, \dots, x_{256})$

linear model: $(w_0, w_1, w_2, \dots, w_{256})$

Features: Extract useful information, e.g.,

intensity and symmetry $\mathbf{x} = (x_0, x_1, x_2)$

linear model: (w_0, w_1, w_2)

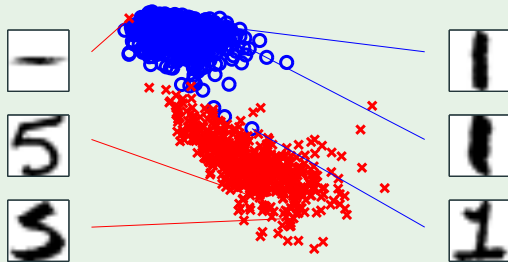


Illustration of features

$$\mathbf{x} = (x_0, x_1, x_2)$$

x_1 : intensity

x_2 : symmetry



A simple learning algorithm - PLA

The perceptron implements

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$$

Given the training set:

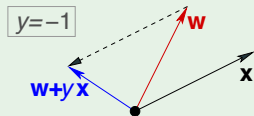
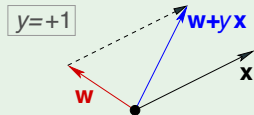
$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$$

pick a **misclassified** point:

$$\text{sign}(\mathbf{w}^T \mathbf{x}_n) \neq y_n$$

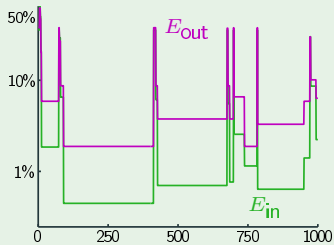
and update the weight vector:

$$\mathbf{w} \leftarrow \mathbf{w} + y_n \mathbf{x}_n$$

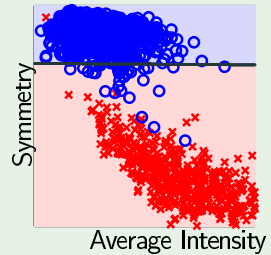


What PLA does

Evolution of E_{in} and E_{out}

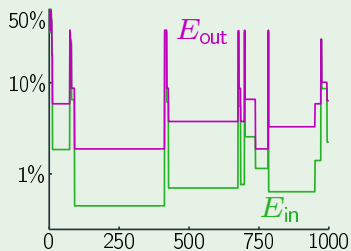


Final perceptron boundary

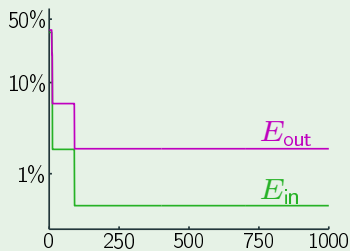


The 'pocket' algorithm

PLA:

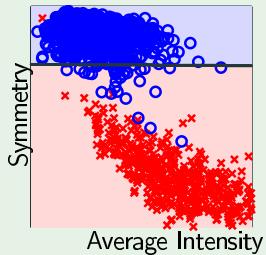


Pocket:



Classification boundary - PLA versus Pocket

PLA:



Pocket:

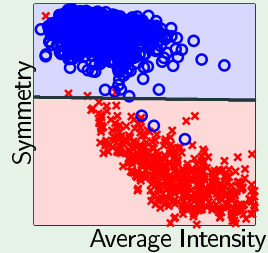


Table of contents

Linear Classification

Linear Regression

Important notions in nonlinear optimization

Back to linear regression

Logistic Regression

Credit again

Classification: Credit approval (yes/no)

Regression: Credit line (dollar amount)

Input: $\mathbf{x} =$

age	23 years
annual salary	\$30,000
years in residence	1 year
years in job	1 year
current debt	\$15,000
...	...

Linear regression output: $h(\mathbf{x}) = \sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x}$

The data set

Credit officers decide on credit lines:

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$$

$y_n \in \mathbb{R}$ is the credit line for customer \mathbf{x}_n .

Linear regression tries to replicate that.

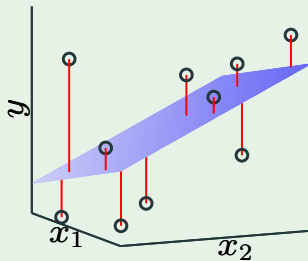
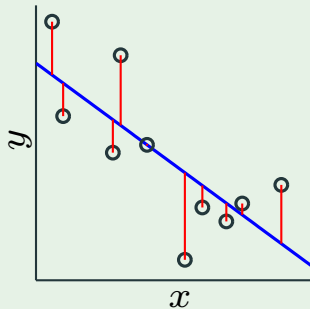
How to measure the error

How well does $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ approximate $f(\mathbf{x})$?

In linear regression, we use squared error $(h(\mathbf{x}) - f(\mathbf{x}))^2$

in-sample error:
$$E_{\text{in}}(h) = \frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n) - y_n)^2$$

Illustration of linear regression



The expression for E_{in}

$$\begin{aligned} E_{\text{in}}(\mathbf{w}) &= \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2 \\ &= \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \end{aligned}$$

where

$$\mathbf{X} = \begin{bmatrix} -\mathbf{x}_1^T \\ -\mathbf{x}_2^T \\ \vdots \\ -\mathbf{x}_N^T \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

The expression for E_{in}

$$\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$$

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2 \quad (1)$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \in \mathbb{R}^N \quad \text{and} \quad \mathbf{X} = \begin{bmatrix} - & \mathbf{x}_1^T & - \\ - & \mathbf{x}_2^T & - \\ & \vdots & \\ - & \mathbf{x}_N^T & - \end{bmatrix} \in \mathbb{R}^{N \times (d+1)}$$

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2, \quad (2)$$

where $\|\cdot\|$ is the Euclidean norm of a vector.

Minimizing E_{in}

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2 \quad (3)$$

$$= \frac{1}{N} \|\mathbf{y} - X\mathbf{w}\|_2^2 \quad (4)$$

$$= \frac{1}{N} (\mathbf{y} - X\mathbf{w})^T (\mathbf{y} - X\mathbf{w}) \quad (\|\mathbf{v}\|_2^2 = \mathbf{v}^T \mathbf{v}) \quad (5)$$

$$= \frac{1}{N} (\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T X^T \mathbf{y} + \mathbf{w}^T X^T X \mathbf{w}). \quad (6)$$

We need to solve the following optimization problem

$$\mathbf{w}_{\text{in}} = \underset{\mathbf{w} \in \mathbb{R}^{d+1}}{\operatorname{argmin}} E_{\text{in}}(\mathbf{w}),$$

where $E_{\text{in}}(\mathbf{w})$ is a single-valued **multivariable** function.

Table of contents

Linear Classification

Linear Regression

Important notions in nonlinear optimization

Back to linear regression

Logistic Regression

Multivariate calculus - Gradient

Consider a function $r(\mathbf{z})$ where \mathbf{z} is a d -vector $\mathbf{z} = [z_1, z_2, \dots, z_d]^T$. The **gradient vector** of this function is given by the **partial derivatives** with respect to each of the independent variables,

$$\nabla r(\mathbf{z}) \equiv g(\mathbf{z}) \equiv \begin{bmatrix} \frac{\partial r}{\partial z_1}(\mathbf{z}) \\ \frac{\partial r}{\partial z_2}(\mathbf{z}) \\ \vdots \\ \frac{\partial r}{\partial z_d}(\mathbf{z}) \end{bmatrix}.$$

Multivariate calculus - Gradient

The gradient is the generalization of the concept of derivative, which captures the local rate of change in the value of a function, in multiple directions.

It is very important to remember that the gradient of a function is only defined if the function is real-valued, that is, if it returns a scalar value.

If the gradient exists at every point, the function is said to be **differentiable**. If each entry of the gradient is continuous, we say the function is **once continuously differentiable**.

Note on notations: $\nabla_{\mathbf{z}} r$ means the gradient of r where the i th partial derivative is taken with respect to z_i . The same notation is used for the directional derivatives defined as $\nabla_{\mathbf{z}} r = \nabla r \cdot \mathbf{z}$.

While the gradient of a function of d variables (i.e. the “first derivative”) is an d -vector, the “second derivative” of an d -variable function is defined by d^2 partial derivatives (the derivatives of the d first partial derivatives with respect to the d variables):

$$\frac{\partial r}{\partial z_i} \left(\frac{\partial r}{\partial z_j} \right) = \frac{\partial^2 r}{\partial z_i \partial z_j}, i \neq j \quad \text{and} \quad \frac{\partial r}{\partial z_i} \left(\frac{\partial r}{\partial z_i} \right) = \frac{\partial^2 r}{\partial^2 z_i}, i = j$$

where $i, j = 1, \dots, d$.

Multivariate calculus - Hessian

If r is single-valued and the partial derivatives $\frac{\partial r}{\partial z_i}$, $\frac{\partial r}{\partial z_j}$ and $\frac{\partial^2 r}{\partial z_i \partial z_j}$ are *continuous*, then $\frac{\partial^2 r}{\partial z_i \partial z_j}$ exists and $\frac{\partial^2 r}{\partial z_i \partial z_j} = \frac{\partial^2 r}{\partial z_j \partial z_i}$.

Therefore the second order partial derivatives can be represented by a square *symmetric* matrix called the **Hessian** matrix:

$$\nabla^2 r(\mathbf{z}) \equiv H(\mathbf{z}) \equiv \begin{pmatrix} \frac{\partial^2 r}{\partial^2 z_1}(\mathbf{z}) & \cdots & \frac{\partial^2 r}{\partial z_1 \partial z_d}(\mathbf{z}) \\ \vdots & & \vdots \\ \frac{\partial^2 r}{\partial z_d \partial z_1}(\mathbf{z}) & \cdots & \frac{\partial^2 r}{\partial^2 z_d}(\mathbf{z}) \end{pmatrix},$$

which contains $d(d+1)/2$ independent elements.

If a function has a Hessian matrix at every point, we say that the function is **twice differentiable**. If each entry of the Hessian is continuous, we say the function is **twice continuously differentiable**.

Multivariate calculus - Hessian

For functions of a vector, the gradient is a vector, and **we cannot take the gradient of a vector**. *Therefore, it is not the case that the Hessian is the gradient of the gradient.* However, this is **almost** true, in the following sense: If we look at the i th entry of the gradient $(\nabla_{\mathbf{z}} r(\mathbf{z}))_i = \frac{\partial r(\mathbf{z})}{\partial z_i}$, and take the gradient with respect to \mathbf{z} , we get

$$\nabla_{\mathbf{z}} \frac{\partial r(\mathbf{z})}{\partial z_i} \equiv \begin{bmatrix} \frac{\partial r}{\partial z_i \partial z_1}(\mathbf{z}) \\ \frac{\partial r}{\partial z_i \partial z_2}(\mathbf{z}) \\ \vdots \\ \frac{\partial r}{\partial z_i \partial z_d}(\mathbf{z}) \end{bmatrix},$$

which is the i th column (or row) of the Hessian. Therefore,

$$\nabla_{\mathbf{z}}^2 = [\nabla_{\mathbf{z}}(\nabla_{\mathbf{z}} r(\mathbf{z}))_1 \quad \nabla_{\mathbf{z}}(\nabla_{\mathbf{z}} r(\mathbf{z}))_2 \quad \cdots \quad \nabla_{\mathbf{z}}(\nabla_{\mathbf{z}} r(\mathbf{z}))_d].$$

Definiteness of a matrix

- $A \in \mathbb{R}^{n \times n}$ is **positive semi-definite**, denoted $A \succeq 0$, if $\mathbf{v}^T A \mathbf{v} \geq 0, \forall \mathbf{v} \in \mathbb{R}_{\geq 0}^n$. If $A = A^T$, all the eigenvalues of A are larger or equal to zero.
- $A \in \mathbb{R}^{n \times n}$ is **positive definite**, denoted $A \succ 0$, if $\mathbf{v}^T A \mathbf{v} > 0, \forall \mathbf{v} \in \mathbb{R}_{> 0}^n$. If $A = A^T$, all the eigenvalues of A are strictly positive.
- $A \in \mathbb{R}^{n \times n}$ is **negative definite**, denoted $A \prec 0$, if $\mathbf{v}^T A \mathbf{v} < 0, \forall \mathbf{v} \in \mathbb{R}_{> 0}^n$. If $A = A^T$, all the eigenvalues of A are strictly negative.
- $A \in \mathbb{R}^{n \times n}$ is **indefinite** if there exists $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}_{> 0}^n$ such that $\mathbf{v}_1^T A \mathbf{v}_1 > 0$ and $\mathbf{v}_2^T A \mathbf{v}_2 < 0$. If $A = A^T$, A has eigenvalues of mixed sign.

Convex functions

A function $r : \mathbb{R}^d \rightarrow \mathbb{R}$ is called **convex** if for any $\mathbf{z}_1, \mathbf{z}_2 \in \mathbb{R}^d$ and every $\lambda \in [0, 1]$, we have

$$f(\lambda \mathbf{z}_1 + (1 - \lambda) \mathbf{z}_2) \leq \lambda f(\mathbf{z}_1) + (1 - \lambda) f(\mathbf{z}_2).$$

- r is **convex** if and only if its Hessian is **positive semidefinite**.
- *For a convex function, any local minimum (optimum) is also a global minimum (optimum).*

Optimality conditions for unconstrained optimization

- Necessary Optimality Conditions

- *First-Order Necessary Conditions*

- If \mathbf{z}^* is a **local minimum** of r and r is once continuously differentiable, **then** $\nabla r(\mathbf{z}^*) = 0$.

- *Second-Order Necessary Conditions*

- If \mathbf{z}^* is a **local minimum** of r and r is twice continuously differentiable, **then** $\nabla^2 r(\mathbf{z}^*) \succeq 0$.

- There may exist points that satisfy these conditions but are not local minima, e.g. $z = 0$ for $r(z) = z^3$, $r(z) = |z|^3$ or $r(z) = -|z|^3$.

- Sufficient Optimality Conditions

- *First-Order Sufficient Conditions*

- If r is once continuously differentiable and convex, and $\nabla r(\mathbf{z}^*) = 0$ **then** \mathbf{z}^* is a **global minimum** of r .

- *Second-Order Sufficient Conditions*

- If r is once continuously differentiable, $\nabla r(\mathbf{z}^*) = 0$ and $\nabla^2 r(\mathbf{z}^*) \succ 0$, **then** \mathbf{z}^* is a **(strict) local minimum** of r .

Table of contents

Linear Classification

Linear Regression

Important notions in nonlinear optimization

Back to linear regression

Logistic Regression

Minimizing E_{in}

Since $E_{\text{in}}(\mathbf{w})$ is differentiable and convex, we can find the global minimum of $E_{\text{in}}(\mathbf{w})$ by requiring $\nabla E_{\text{in}}(\mathbf{w}) = 0$. We have

$$\begin{aligned}\nabla_{\mathbf{w}} \left(\frac{1}{N} (\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T X^T \mathbf{y} + \mathbf{w}^T X^T X \mathbf{w}) \right) &= 0 \\ \iff \frac{2}{N} (X^T X \mathbf{w} - X^T \mathbf{y}) &= 0 \iff X^T X \mathbf{w} = X^T \mathbf{y} \\ \iff \mathbf{w} = (X^T X)^{-1} X^T \mathbf{y} &\quad (\text{assuming } X^T X \text{ is invertible})\end{aligned}$$

where we used the following gradient identities

$$\nabla_{\mathbf{w}}(\mathbf{w}^T A \mathbf{w}) = (A + A^T) \mathbf{w}, \quad \nabla_{\mathbf{w}}(\mathbf{w}^T \mathbf{b}) = \mathbf{b}.$$

The Hessian of $E_{\text{in}}(\mathbf{w})$ is given by

$$H(\mathbf{w}) = \nabla_{\mathbf{w}} \left(\frac{2}{N} (X^T X \mathbf{w} - X^T \mathbf{y}) \right) = \frac{2}{N} X^T X, \quad (7)$$

where we used $\nabla_{\mathbf{w}}(A \mathbf{w}) = A$. $H(\mathbf{w})$ is positive semi-definite.

Minimizing E_{in}

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

$$\nabla E_{\text{in}}(\mathbf{w}) = \frac{2}{N} \mathbf{X}^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) = \mathbf{0}$$

$$\mathbf{X}^\top \mathbf{X} \mathbf{w} = \mathbf{X}^\top \mathbf{y}$$

$$\mathbf{w} = \mathbf{X}^\dagger \mathbf{y} \quad \text{where} \quad \mathbf{X}^\dagger = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

\mathbf{X}^\dagger is the 'pseudo-inverse' of \mathbf{X}

The linear regression algorithm

- 1: Construct the matrix X and the vector \mathbf{y} from the data set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ as follows

$$\underbrace{X = \begin{bmatrix} -\mathbf{x}_1^\top - \\ -\mathbf{x}_2^\top - \\ \vdots \\ -\mathbf{x}_N^\top - \end{bmatrix}}_{\text{input data matrix}}, \quad \underbrace{\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\text{target vector}}.$$

- 2: Compute the pseudo-inverse $X^\dagger = (X^\top X)^{-1} X^\top$.
- 3: Return $\mathbf{w} = X^\dagger \mathbf{y}$.

Learning curves in linear regression

Expected E_{out} and E_{in}

Data set \mathcal{D} of size N

Expected out-of-sample error $\mathbb{E}_{\mathcal{D}}[E_{\text{out}}(g^{(\mathcal{D})})]$

Expected in-sample error $\mathbb{E}_{\mathcal{D}}[E_{\text{in}}(g^{(\mathcal{D})})]$

How do they vary with N ?

Learning curves in linear regression

Linear regression case

Noisy target $y = \mathbf{w}^{*T} \mathbf{x} + \text{noise}$

Data set $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$

Linear regression solution: $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

In-sample error vector = $\mathbf{X} \mathbf{w} - \mathbf{y}$

'Out-of-sample' error vector = $\mathbf{X} \mathbf{w} - \mathbf{y}'$

Learning curves in linear regression

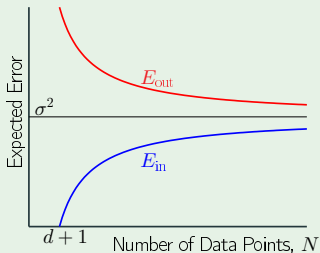
Learning curves for linear regression

Best approximation error = σ^2

Expected in-sample error = $\sigma^2 \left(1 - \frac{d+1}{N}\right)$

Expected out-of-sample error = $\sigma^2 \left(1 + \frac{d+1}{N}\right)$

Expected generalization error = $2\sigma^2 \left(\frac{d+1}{N}\right)$



Linear regression for classification

Linear regression learns a real-valued function $y = f(\mathbf{x}) \in \mathbb{R}$

Binary-valued functions are also real-valued! $\pm 1 \in \mathbb{R}$

Use linear regression to get \mathbf{w} where $\mathbf{w}^T \mathbf{x}_n \approx y_n = \pm 1$

In this case, $\text{sign}(\mathbf{w}^T \mathbf{x}_n)$ is likely to agree with $y_n = \pm 1$

Good initial weights for classification

Linear regression boundary

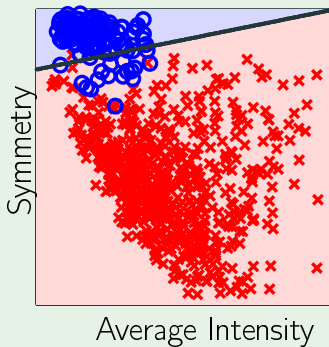


Table of contents

Linear Classification

Linear Regression

Important notions in nonlinear optimization

Back to linear regression

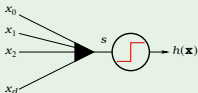
Logistic Regression

A third linear model

$$s = \sum_{i=0}^d w_i x_i$$

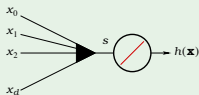
linear classification

$$h(\mathbf{x}) = \text{sign}(s)$$



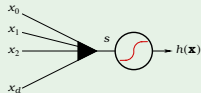
linear regression

$$h(\mathbf{x}) = s$$



logistic regression

$$h(\mathbf{x}) = \theta(s)$$



The logistic function θ

The formula:

$$\theta(s) = \frac{e^s}{1 + e^s}$$



soft threshold: uncertainty

sigmoid: flattened out 's'

Probability interpretation

$h(\mathbf{x}) = \theta(s)$ is interpreted as a probability

Example. Prediction of heart attacks

Input \mathbf{x} : cholesterol level, age, weight, etc.

$\theta(s)$: probability of a heart attack

The signal $s = \mathbf{w}^T \mathbf{x}$ "risk score"

Genuine probability

Data (\mathbf{x}, y) with **binary** y , generated by a noisy target:

$$P(y \mid \mathbf{x}) = \begin{cases} f(\mathbf{x}) & \text{for } y = +1; \\ 1 - f(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

The target $f : \mathbb{R}^d \rightarrow [0, 1]$ is the probability

Learn $g(\mathbf{x}) = \theta(\mathbf{w}^\top \mathbf{x}) \approx f(\mathbf{x})$

Genuine probability

Data (\mathbf{x}, y) with **binary** y , generated by a noisy target:

$$P(y \mid \mathbf{x}) = \begin{cases} f(\mathbf{x}) & \text{for } y = +1; \\ 1 - f(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

The target $f : \mathbb{R}^d \rightarrow [0, 1]$ is the probability

Learn $g(\mathbf{x}) = \theta(\mathbf{w}^\top \mathbf{x}) \approx f(\mathbf{x})$

The data does not give us the value of f explicitly. It gives us samples generated by this probability. How do we learn from such data?

Error measure

For each (\mathbf{x}, y) , y is generated by probability $f(\mathbf{x})$

Plausible error measure based on **likelihood**:

If $h = f$, how likely to get y from \mathbf{x} ?

$$P(y | \mathbf{x}) = \begin{cases} h(\mathbf{x}) & \text{for } y = +1; \\ 1 - h(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

Formula for likelihood

Since the data points are assumed to be (conditionally) independently generated, the probability of observing all the y_n 's in the data set from the corresponding \mathbf{x}_n is given by

$$\prod_{n=1}^N P(y_n|\mathbf{x}_n).$$

where

$$P(y|\mathbf{x}) = \begin{cases} h(\mathbf{x}) & \text{for } y = +1; \\ 1 - h(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

The method of *maximum likelihood* selects the hypothesis h which maximizes this probability.

From likelihood to E_{in}

$$\begin{aligned}\text{Maximize } \prod_{n=1}^N P(y_n | \mathbf{x}_n) &\equiv \text{Maximize } \ln \left(\prod_{n=1}^N P(y_n | \mathbf{x}_n) \right) \\ &\equiv \text{Maximize } \frac{1}{N} \ln \left(\prod_{n=1}^N P(y_n | \mathbf{x}_n) \right) \\ &\equiv \text{Minimize } -\frac{1}{N} \ln \left(\prod_{n=1}^N P(y_n | \mathbf{x}_n) \right)\end{aligned}$$

Furthermore, we can write

$$\begin{aligned}& -\frac{1}{N} \ln \left(\prod_{n=1}^N P(y_n | \mathbf{x}_n) \right) \\ &= \frac{1}{N} \sum_{n=1}^N \ln \left(\frac{1}{P(y_n | \mathbf{x}_n)} \right) \\ &= \frac{1}{N} \sum_{n=1}^N \mathbb{1}\{y_n = +1\} \ln \left(\frac{1}{h(\mathbf{x}_n)} \right) + \mathbb{1}\{y_n = -1\} \ln \left(\frac{1}{1 - h(\mathbf{x}_n)} \right)\end{aligned}$$

From likelihood to E_{in}

We have $h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$, where $\theta(s) = \frac{e^s}{1+e^s} = \frac{1}{1+e^{-s}}$ with $\theta(-s) = 1 - \theta(s)$. So, we can write

$$\begin{aligned} & \frac{1}{N} \sum_{n=1}^N \mathbb{1}\{y_n = +1\} \ln \left(\frac{1}{h(\mathbf{x}_n)} \right) + \mathbb{1}\{y_n = -1\} \ln \left(\frac{1}{1 - h(\mathbf{x}_n)} \right) \\ &= \frac{1}{N} \sum_{n=1}^N \mathbb{1}\{y_n = +1\} \ln \left(\frac{1}{\theta(\mathbf{w}^T \mathbf{x}_n)} \right) + \mathbb{1}\{y_n = -1\} \ln \left(\frac{1}{1 - \theta(\mathbf{w}^T \mathbf{x}_n)} \right) \\ &= \frac{1}{N} \sum_{n=1}^N \mathbb{1}\{y_n = +1\} \ln \left(\frac{1}{\theta(\mathbf{w}^T \mathbf{x}_n)} \right) + \mathbb{1}\{y_n = -1\} \ln \left(\frac{1}{\theta(-\mathbf{w}^T \mathbf{x}_n)} \right) \\ &= \frac{1}{N} \sum_{n=1}^N \ln \left(\frac{1}{\theta(y_n \mathbf{w}^T \mathbf{x}_n)} \right) \quad (\text{i.e. } P(y_n | \mathbf{x}_n) = \theta(y_n \mathbf{w}^T \mathbf{x}_n)) \\ &= \underbrace{\frac{1}{N} \sum_{n=1}^N \ln \left(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n} \right)}_{E_{\text{in}}(\mathbf{w})} \end{aligned}$$

Cross-entropy

For two probability distributions with binary outcomes $\{p, 1 - p\}$ and $\{q, 1 - q\}$, the cross-entropy (from information theory) is

$$p \log \frac{1}{q} + (1 - p) \log \frac{1}{1 - q}.$$

The in-sample error above corresponds to a cross-entropy error measure on the data point (\mathbf{x}_n, y_n) , with $p = \mathbb{1}\{y_n = +1\}$ and $q = h(\mathbf{x}_n)$.

Formula for likelihood

$$P(y \mid \mathbf{x}) = \begin{cases} h(\mathbf{x}) & \text{for } y = +1; \\ 1 - h(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

Substitute $h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$, noting $\theta(-s) = 1 - \theta(s)$



$$P(y \mid \mathbf{x}) = \theta(y \mathbf{w}^T \mathbf{x})$$

Likelihood of $\mathcal{D} = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ is

$$\prod_{n=1}^N P(y_n \mid \mathbf{x}_n) = \prod_{n=1}^N \theta(y_n \mathbf{w}^T \mathbf{x}_n)$$

Maximizing the likelihood

Minimize
$$-\frac{1}{N} \ln \left(\prod_{n=1}^N \theta(y_n \mathbf{w}^\top \mathbf{x}_n) \right)$$

$$= \frac{1}{N} \sum_{n=1}^N \ln \left(\frac{1}{\theta(y_n \mathbf{w}^\top \mathbf{x}_n)} \right) \quad \left[\theta(s) = \frac{1}{1 + e^{-s}} \right]$$

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \underbrace{\ln \left(1 + e^{-y_n \mathbf{w}^\top \mathbf{x}_n} \right)}_{e(h(\mathbf{x}_n), y_n)} \quad \text{"cross-entropy" error}$$

How to minimize E_{in}

For logistic regression,

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln \left(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n} \right) \quad \leftarrow \text{iterative solution}$$

Compare to linear regression:

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2 \quad \leftarrow \text{closed-form solution}$$

Summary of Linear Models

