binary

Previously in classification, we measured error by

$$MAE = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}_{\hat{y}_i \neq y_i}$$  AKA  misclassification error

Does this work for both in-sample & oos?  YES  → "negative"
→ "positive"

But this hides a lot of what's going on! There are two types of errors. (1) Saying a the 0 is a 1 (false positive) and (2) " " " the 1 is a 0 (false negative).

How can we visualize this?  In a 2×2 table:



prediction ($\hat{y}$)

|       |   | 0 | 1 |     |
|-------|---|----|----|-----|
| truth | 0 | TN | FP | #N |
| (y)   | 1 | FN | TP | #P |
|       |   | #PN | #PP | n |

predicted   predicted
negative    positive

$$\text{misclassification error} = \frac{FP + FN}{n}$$

$$\text{Accuracy} = 1 - \text{misclassification error}$$
$$= \frac{TP + TN}{n}$$

there are many metrics!

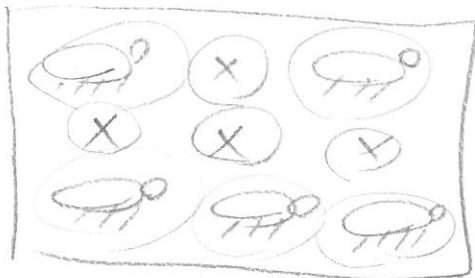$$\text{Precision} = \frac{TP}{\#PP}$$  the prop. of those you predicted positive actually positive?

$$\text{Recall / Sensitivity} = \frac{TP}{\#P}$$  " " " those truly positive were predicted positive?

Example:



$$Precision = \frac{2}{2} = 100\%$$

$$Recall = \frac{2}{5} = 90\%$$



$$precision = \frac{5}{9} = 56\%$$

$$recall = \frac{5}{5} = 100\%$$

Both situations are bad!  How to combine both together?

$$F_1 score := \frac{2}{\frac{1}{Recall} + \frac{1}{precision}} \overset{e.g}{=} \frac{2}{\frac{1}{1} + \frac{1}{.56}} = 72\% \quad \text{(popular,} \\ \text{s.cc. also popular)}$$

More:

false discovery rate

$$FDR = 1 - precision = \frac{FP}{\#PP}$$

false omission rate

$$FOR = \frac{FP}{\#PN}$$

In-sample and OOS

Are these confusion tables for classification with $K > 2$ levels?

Yes!! DEMO

In my experience more important. They operate on the columns of the confusion table. This means FPR tells you if $\hat{y} = 1$, whis your error rate & FOR ′′′′′ $\hat{y} = 1$, ′′′′′′′

DEMO

Credit Loan question

① How much will person pay back?

② Will they pay back dtly?

③ What is the chance they pay back?

Type of Supervised Learning Task

Regression

Classification

Probabilistic classification

Probabilistic Classification : the response is still the label. Let's focus on binary

Precisely...

$$y = \underbrace{t(z_1, \ldots, z_t)}_{\in \{0,1\}} = \underbrace{f(x_1, \ldots, x_p)}_{\in \{0,1\}} + \underbrace{\delta}_{\in \{+1, -1\}} = \underbrace{h^u(x_1, \ldots, x_p)}_{\in \{0,1\}} + \underbrace{\varepsilon}_{\in \{+1, -1\}} = \underbrace{g(x_1, \ldots, x_p)}_{\in \{0,1\}} + \underbrace{e}_{\in \{+1, -1\}}$$

$f, h^u, g$ all have $y = \{0,1\}$ as their output. They do __not__ learn probs.

Consider an alternative construction. Let $Y$ represent the r.v. whose realization is $y$. Instead of the above

$$Y \sim \text{bernoulli}\left( f_{pr}(x_1, \ldots, x_p) \right)$$

prob this $\vec{x}$ will be 1

for a given $\vec{x}_i = [x_{i1}, \ldots, x_{ip}]$, the above r.v. is realized to $y_i = 1$ or $y_i = 0$. Same for

$$Y \sim \text{bernoulli}\left( h^u_{pr}(x_1, \ldots, x_p) \right) \quad \& \quad Y \sim \text{bernoulli}\left( g_{pr}(x_1, \ldots, x_p) \right)$$

Where are the error terms: $\delta, \varepsilon, e$? Not needed in this construction. The goal is to model the prob's. So errors now are differences with the true prob. function $t(z_1, \ldots, z_t) = t_{pr}(z_1, \ldots, z_t) \leftarrow$ why?

$$Y = Bern \left( t_{pr} (z_1, \ldots, z_t) \right) \iff y = t(z_1, \ldots, z_t)$$

No error means prob's are 0 or 1!

What is the difference between $t_{pr}(z_1, \ldots, z_t)$ & $f_{pr}(x_1, \ldots, x_p)$?

$f_{pr} \approx t$ (hopefully) but $f_{pr}$ will return some values $\in (0,1)$ since the error due to ignorance will be captured as a non 0 or 1 prob!

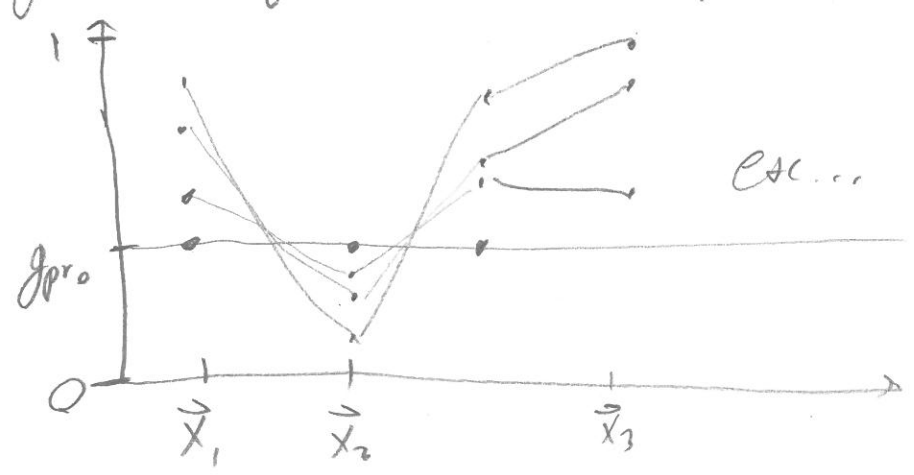e.g. $t(\vec{z}) = 1$ and $g_{pr}(\vec{x}) = 0.8$

these correspond

this means $g(\vec{x})$ will be correct most of the time. Only $\frac{1}{5}$ on avg. will $\hat{y} = 0$ when $y = 1$.

What is the difference between $f_{pr}, h^*_{pr}, g_{pr}$ ?

the values of these functions are further away from perfect 0's & 1's.
Null model
$$g_{pr,0} = \frac{1}{h} \Sigma y_i = \hat{p}$$

e.g. let's say there are even # of 0's and 1's in $\mathbb{D}$ with $n = 200$.



etc...

$f_{pr}$ is closer to 0's, 1's then $h^*_{pr}$ which is closer to 0's, 1's then $g_{pr}$.

As your model gets worse and worse the prob. estimates move further from 0's & 1's to values closer to $\hat{y}_{pr,0}$, the overall avg.

OK... how do we create $f_{pr}$. We need an algorithm R. First let's do some more prob. At the best we can know...
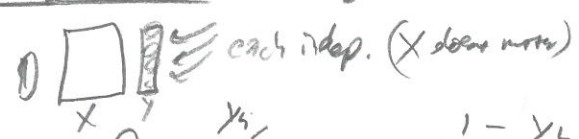
$$Y_1 \sim Bern\left(f_{pr}(\vec{x}_1)\right), \quad Y_2 \sim Bern\left(f_{pr}(\vec{x}_2)\right), \ldots, Y_n \sim Bern\left(f_{pr}(\vec{x}_n)\right)$$

$$\underset{\parallel}{} \qquad \underset{\parallel}{}$$

$$\left(f_{pr}(\vec{x}_1)\right)^{y_1}\left(1-f_{pr}(\vec{x}_1)\right)^{1-y_1} \qquad \cdots \qquad \cdots$$

What is $P\left(Y_1, \ldots, Y_n\right)$ is the "joint mass function"?

Unknown unless we know the <u>dependence</u> structure between $Y_1, \ldots, Y_n$.

Typically a big <u>assumption is made</u> ... <u>independence</u> which gives us:

each idep. (X doesn't matter)

$$P\left(Y_1, \ldots, Y_n\right) = f_{pr}(\vec{x}_1)^{y_1}\left(1-f_{pr}(\vec{x}_1)\right)^{1-y_1} \cdot \ldots \cdot f_{pr}(\vec{x}_n)^{y_n}\left(1-f_{pr}(\vec{x}_n)\right)^{1-y_n}$$

$$= \prod_{i=1}^{n} f_{pr}(\vec{x}_i)^{y_i}\left(1-f_{pr}(\vec{x}_i)\right)^{1-y_i}$$

Now... if we are trying to estimate $f_{pr}$, what are the "knowns"? The $y_1, y_2, \ldots, y_n$. So our goal now is to come up with $f_{pr}$ s.t. the prob. $P(Y_1, \ldots, Y_n)$ (AKA the "likelihood") is <u>MAXIMIZED</u>.

But of course, $f_{pr}$ is arbitrarily complicated with interactions & nonlinearities. So let's make an assumption on candidate models.

$\mathcal{H} = \{$ set of all candidate prob. functions $\}$ of which $h^*_{pr}$ is the closest to $f_{pr}$. Then we use an alg. $\mathcal{A}$ to pick $g_{pr}$ which is the best and hope it's close to $h^*_{pr}$.

What can we use for $\mathcal{H}$? How about $\mathcal{H} = \{\vec{w} \cdot \vec{x} : \vec{w} \in \mathbb{R}^{P+1}\}$?

Why doesn't this work? $\vec{w} \cdot \vec{x} \in \mathbb{R}$ and prob's $\in [0,1]$

What about $\mathcal{H} = \{\mathbb{1}_{\vec{w} \cdot \vec{x} \geq 0}\}$? Only 0 or 1... nothing in between!

What if we wanted to keep the <u>linear</u> model but wanted $\underset{\text{output}}{\text{smooth}} \in (0,1)$?
↳ generalized linear model (GLM)

we need a "link function" $\phi(\vec{w} \cdot \vec{x})$ whose range is $(0,1)$ monotonically & smoothly.

$\left(\text{i.e. if } \vec{w} \cdot \vec{x} \uparrow \Rightarrow \text{prob est.} \uparrow\right)$ Also, $\phi(\vec{w} \cdot \vec{x}) \neq 0 \text{ or } 1$ Why?

we can never be sure!! There is information we don't know!!

There are many possible $\phi$ functions! (These are also called Activation functions in neural nets.)

The most popular is the logistic function:

$$\phi(u) = \frac{e^u}{1+e^u} = \frac{1}{1+e^{-u}}$$

Also the probit function $\phi(u) = \Phi^{-1}(u)$ — inverse CDF of std. normal

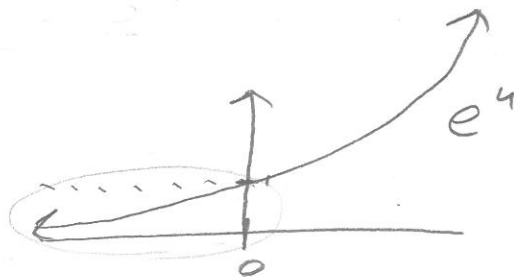or the "complementary log-log":

$$\phi(u) = 1 - e^{-e^u} =$$

Why is this $\in (0,1)$?

$$e^u \in (0,\infty)$$
$$-e^u \in (-\infty, 0)$$
$$e^{-e^u} \in (0,1)$$
$$1 - e^{-e^u} \in (0,1)$$



or the hyperbolic tangent:

$$\phi(u) = \tanh(u) := \frac{e^u - e^{-u}}{e^u + e^{-u}}$$

So let's use logistic $\phi$:

$$\mathcal{H} = \left\{ \frac{e^{\vec{w}\cdot\vec{x}}}{1 + e^{\vec{w}\cdot\vec{x}}} \; : \; \vec{w} \in \mathbb{R}^{p+1} \right\}$$

What does $A$ do? It maximizes the likelihood:

$$\vec{b} := \operatorname*{argmax}_{\vec{w}} \left\{ P(Y_1, \dots, Y_n) \right\} = \operatorname*{argmax}_{\vec{w}} \left\{ \prod_{i=1}^{n} \left( \frac{e^{\vec{w}\cdot\vec{x}_i}}{1 + e^{\vec{w}\cdot\vec{x}_i}} \right)^{Y_i} \left( 1 - \frac{e^{\vec{w}\cdot\vec{x}_i}}{1 + e^{\vec{w}\cdot\vec{x}_i}} \right)^{1 - Y_i} \right\}$$

Since we are seeking the argmax of some function $\alpha(\vec{w})$ we can easily just find the argmax of $\ln(\alpha(\vec{w}))$ as well:

$$= \operatorname*{argmax}_{\vec{w}} \left\{ \prod_{i=1}^{n} \left( \frac{1}{1+e^{-\vec{w}\cdot\vec{x}_i}} \right)^{y_i} \left( \frac{1}{1+e^{\vec{w}\cdot\vec{x}_i}} \right)^{1-y_i} \right\} = \operatorname*{argmax}_{\vec{w}} \left\{ \prod_{i=1}^{n} \frac{1}{1+e^{\vec{w}\cdot\vec{x}_i}} \prod_{i=1}^{n} \left( \frac{1+e^{\vec{w}\cdot\vec{x}_i}}{1+e^{-\vec{w}\cdot\vec{x}_i}} \right)^{y_i} \right\}$$

$$\left( \frac{1}{1+e^{\vec{w}\cdot\vec{x}_i}} \right) \left( 1+e^{\vec{w}\cdot\vec{x}_i} \right)^{y_i}$$

$$= \begin{cases} \left( 1+e^{-\vec{w}\cdot\vec{x}_i} \right)^{-1}_i & \text{if} \quad y_i = 1 \\[2em] \left( 1+e^{\vec{w}\cdot\vec{x}_i} \right)^{-1} & \text{if} \quad y_i = 0 \end{cases}$$

**p91
LFD**

Consequence of using
logistic link!

$$= \left( 1 + e^{(1-2y_i)\vec{w}\cdot\vec{x}} \right)^{-1}$$

$$= -(1-2y_i)$$

let $z_i = 2y_i - 1$

$y_i = 0 \Rightarrow z_i = -1$
$y_i = 1 \Rightarrow z_i = +1$

$$= \operatorname*{argmax}_{\vec{w}} \left\{ \prod_{i=1}^{n} \left( 1 + e^{-z_i \vec{w}\cdot\vec{x}_i} \right)^{-1} \right\}$$

Note: if we are taking $\operatorname{argmax} \left\{ v(t) \right\} = \operatorname{argmax} \left\{ \ln(v(t)) \right\}$
since $\ln$ is a monotonic increasing transformation

$$= \operatorname*{argmax}_{\vec{w}} \left\{ \ln\left( \quad \right) \right\} = \operatorname*{argmax}_{\vec{w}} \left\{ -\sum_{i=1}^{n} \ln\left( 1 + e^{-z_i \vec{w}\cdot\vec{x}_i} \right) \right\}$$

$$= \operatorname*{argmin}_{\vec{w}} \left\{ \sum_{i=1}^{n} \ln\left( 1 + e^{-z_i \vec{w}\cdot\vec{x}_i} \right) \right\}$$

Now we can take $\frac{d}{d\vec{w}} \left[ \quad \right] \overset{\text{set}}{=} 0$ to solve for $\vec{b}$.