Classification / probability estimates.

↑
want {0,1}      ↑ want $0 \leq p \leq 1$    linear regression often
   {0,... k}                                 does not
answer                                        a good fit           [crossed out]

## Logistic regression    Example.

Logistic function   $\dfrac{e^x}{1+e^x}$



$0 \leq \quad \leq 1$

die ┬ . . . . . .
live ├ x x x . . . → age

↑ linear model ẑ:
$y = mx + c$

↑ fit a logistic function instead

if the probability of an event is $p$
the odds of the event are $\dfrac{p}{1-p}$, so  $\dfrac{p(x)}{1-p(x)} = e^{mx+c}$.

$$p(x) = \dfrac{e^{mx+c}}{1+e^{mx+c}}.$$

Logistic function ↔ log(odds) is linear.
              logit

## Estimating the coefficients  $mx+c$ / $\hat{\beta_0} + \hat{\beta_1} a$.

instead of least squares error use  <u>maximum likelihood</u>.

likelihood function   $\ell(m,c) = \prod\limits_{i, y_i = 1} p(x_i) \prod\limits_{i, y_i = 0} (1-p(x_i))$.

intuition : $\ell = 1$ best case, all predictions correct
         $\ell = 0$ at least one prediction completely opposite.

<u>Fact</u> : in many cases there are fast algorithms to find $m, c$.

end up with  $p(x) = \dfrac{e^{mx+c}}{1+e^{mx+c}}$ } estimates prob of event occurring

## Multiple Logistic regression (multiple dim input, one prob output).

use log odds :  $\log\left(\dfrac{p(x)}{1-p(x)}\right) = \beta_0 + \beta_1 X_1 + \cdots + \beta_d X_d.$

equivalent to: $p(x) = \dfrac{e^{\beta_0 + \beta_1 X_1 + \cdots + \beta_d X_d}}{1 + e^{\beta_0 + \beta_1 X_1 + \cdots + \beta_d X_d}}$

__Multinomial logistic regression__ (need to predict one of K things... ) .

$\{1, \ldots, K\}$ .

choose a __base class__ K.

then for $1 \le k \le K-1$ use: $P(Y=k \mid X=x) = \dfrac{e^{\beta_{k0} + \beta_{k1} x_1 + \cdots + \beta_{kd} x_d}}{1 + \sum\limits_{\ell=1}^{K-1} e^{\beta_{\ell 0} + \beta_{\ell 1} x_1 + \cdots + \beta_{\ell d} x_d}}$

for K: $P(Y=k \mid X=x) = \dfrac{1}{1 + \sum\limits_{\ell=1}^{K-1} e^{\beta_{\ell 0} + \beta_{\ell 1} x_1 + \cdots + \beta_{\ell d} x_d}}$ .

then.

$\log\left( \dfrac{P(Y=k \mid X=x)}{P(Y=K \mid X=x)} \right) = \beta_{k0} + \beta_{k1} x_1 + \cdots + \beta_{kd} x_d$ . $\leftarrow$ i.e. _log_ conditional prob are _linear_ .

__Fact__ doesn't matter which class you choose to be base class

but coeffs $\beta_{ij}$ depend on choice of base class.

__Alternative__ softmax coding (no base class — symmetric) .

$P(Y=k \mid X=x) = \dfrac{e^{\beta_{k0} + \beta_{k1} x_1 + \cdots + \beta_{kd} x_d}}{\sum\limits_{\ell=1}^{K} e^{\beta_{\ell 0} + \beta_{\ell 1} x_1 + \cdots + \beta_{\ell d} x_d}}$ .

log odds ratio is: $\log\left( \dfrac{P(Y=k \mid X=x)}{P(Y=k' \mid X=x)} \right) = (\beta_{k0} - \beta_{k'0}) + (\beta_{k1} - \beta_{k'1}) x_1 + \cdots + (\beta_{kd} - \beta_{k'd}) x_d$ .

# Generalised linear models

use $X$ to predict $Y$ ← quantitative, linear regression

← qualitative, logistic regression

space $Y$ is a count (e.g. # of people).

## Example   Poisson regression

recall   Poisson dist : $\mathbb{P}(X=n) = \dfrac{e^{-\lambda}\lambda^n}{n!}$   for $n \in \mathbb{N}_0$.   $\mathbb{E}(X) = \lambda$.

$Var(X) = \lambda$.

model count as a Poisson dist with parameter $\lambda$.

when $\lambda(x)$   usual choice   $\lambda = e^{\beta_0 + \beta_1 x_1 + \cdots + \beta_d x_d}$.

$\iff \log(\lambda) = \beta_0 + \beta_1 x_1 + \cdots + \beta_d x_d$

i.e. choose $\log(\lambda)$ to be a linear function of the ~~data~~ data space.

to find $\beta_i$ use maximum likelihood.

$$\ell(\beta_0, \beta_1, \ldots, \beta_d) = \prod_{i=1}^{n} \frac{e^{-\lambda(x_i)} \lambda(x_i)^{y_i}}{(y_i)!}$$

i.e. find $\beta_i$ that maximise this ↑ this is the probability that you obtained that data set from the given value of $\lambda$.

Fact the $\beta_i$ can be estimated efficiently

- interpretation of coefficients: $\beta_i$ big means column more important.
  increase $x_i \mapsto x_i + 1$ then $\mathbb{E}(Y)$ increases by a factor of $e^{\beta_i}$.

- Poisson model has $\lambda = E(X) = Var(X)$ often better fit than underlying
assumption of linear regression model where var is constant.
- positive integer output, discrete, no negative values.

# Generalised linear models

regression models:   linear → quantities
logistic → categorical
Poisson → $\mathbb{N}_0$.

common features:

- $f(x)$ predicts $Y$.

- model analysis usually assumes $Y$ has
  - Gaussian dist for linear
  - Bernoulli $q_{p_i}$ for categorical
  - Poisson for Poisson

- linear: $\mathbb{E}(Y|x) = \beta_0 + \beta_1 x_1 + \cdots \beta_d x_d$.

  logistic: $\mathbb{E}(Y|x) = \mathbb{P}(Y=1|x)$.

  $$= \frac{e^{\beta_0 + \beta_1 x_1 + \cdots + \beta_d x_d}}{1 + e^{\beta_0 + \beta_1 x_1 + \cdots + \beta_d x_d}} \quad \Longleftrightarrow$$

  equivalent log odds are
  
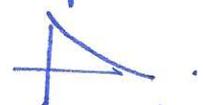  $\log\left(\frac{\mathbb{E}(Y|x)}{1-\mathbb{E}(Y|x)}\right) = \beta_0 + \beta_1 x_1 + \cdots + \beta_d x_d$

  Poisson $\mathbb{E}(Y|x) = \lambda(x) = e^{\beta_0 + \beta_1 x_1 + \cdots + \beta_d x_d} \quad \Longleftrightarrow$

  $\log\left(\underbrace{\mathbb{E}(Y|x)}_{\lambda(x)}\right) = \beta_0 + \beta_1 x_1 + \cdots + \beta_d x_d$.

__notation__  link function transforms the expected
mean so that it is linear for:

linear $\eta(\mu) = \mu$    logistic $\eta(\mu) = \log\left(\frac{\mu}{1-\mu}\right)$    Poisson $\eta(\mu) = \log(\mu)$.

__Fact__  these dists are all examples of exponential family dists.

other examples    • exponential dist

$X_i$ indep exp $\beta$
then $\sum X_i \sim \Gamma(n, \beta)$  $d \in \mathbb{N}$
$\Gamma(1, \beta) = \text{Exp}(\beta)$.  get
$\Gamma(k/2, 1/2) \text{ is } \chi^2(k)$  Poisson..

• Gamma dist  $\Gamma(\alpha, \theta)$  $f(x) = \frac{1}{\Gamma(\alpha)\theta^\alpha} x^{\alpha-1} e^{-x/\theta}$
$\alpha$ shape
$\theta$ scale.
mean $\alpha\theta$
var $\alpha\theta^2$

$\alpha < \theta$.
$\alpha > \theta$

• negative binomial dist

these are regression schemes for all them w/ different choice of $\eta$.

→ discrete analog of $\Gamma$.  $k \mapsto \binom{k+r-1}{k}(1-p)^k p^r$
NB$(r, p)$ mean $\frac{r(1-p)}{p}$   var $\frac{r(1-p)}{p^2}$

Q: how many times
do I need to run my
Bernoulli press before I set a works?   NB$(r, p)$
#of times it occurs.  prob of Bernoulli trys