

Classification

Training data $\{(x_i, y_i)\}_{i=1}^N$

↑
predictors,
features

↑
response,
labels

y_i : categorical/
qualitative

x_{ij} : continuous,
ordinal,
categorical

Goal: Learn from training data a good classification
function / classifier $\hat{y} = \hat{f}(x)$ or $\hat{y} = \hat{C}(x)$

Ex: $y_i \in C = \left\{ \underbrace{\text{"good email"}}_{\text{"0"}}, \underbrace{\text{"spam"}}_{\text{"1"}} \right\}$

$\tilde{x}_i = [x_{i1}, \dots, x_{ip}]^T$

↑
total word
count

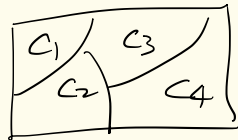
↑
number of irregular
spelling

Step 1: Use a naive Bayes to learn $\hat{y} = \hat{C}_{Tr}(x)$

Step 2: For incoming data x^o , predict if spam or not using $\hat{C}_{Tr}(x^o)$.

Error metric for multiclass classifiers:

• error rate : $\frac{1}{N} \sum_{i=1}^{N+M} \mathbb{1}(y_i \neq \hat{y}_i)$



Bayesian Decision Theory (Murphy 5.7)

Loss Function $L(y, c)$, e.g., $L = \mathbb{1}[y \neq c]$ 0-1 loss

\uparrow true label
 \uparrow label/class (to be) assessed,
label "action"

Expected loss $p(c|\underline{x}) = \mathbb{E}[L(y, c) | \underline{x}]$

$$= \sum_{y \in \{c_1, \dots, c_k\}} L(y, c) p(y|\underline{x})$$

Bayes Decision Rule:

$$\hat{C}(\underline{x}) = \arg \min_c p(c|\underline{x})$$

$$L(y, c) = \mathbb{1}[y \neq c]$$

$$= \arg \min_c \sum_y L(y, c) p(y|\underline{x})$$

$$= \arg \min_c (1 - p(c|\underline{x}))$$

$$= \arg \max_{y \in \{c_1, \dots, c_k\}} p(y|\underline{x})$$

Maximum a posteriori (MAP) estimator

Will be used for "logistic regression" estimation.

Error Metrics for binary classifiers:

- False positive rate (FPR, Type-1 error)
"alarm"
- False negative rate (FNR, Type-2 error)
"miss"

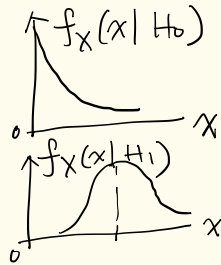
| | | |
|----------|--------------|-------|
| | Ground truth | |
| | H_0 | H_1 |
| Decision | 0 | FN |
| | 1 | TP |

[TPR = sensitivity, recall. TNR = specificity, selectivity]

Ex: (Simple thresholding rule for COVID detection)

$$H_0 \text{ (not infected): } f_X(x|H_0) = \lambda \cdot e^{-\lambda x}, \quad x \geq 0$$

$$H_1 \text{ (infected): } f_X(x|H_1) = \frac{1}{\sqrt{2\pi}\sigma_1^2} e^{-\frac{(x-\mu)^2}{2\sigma_1^2}}, \quad x \in \mathbb{R}$$

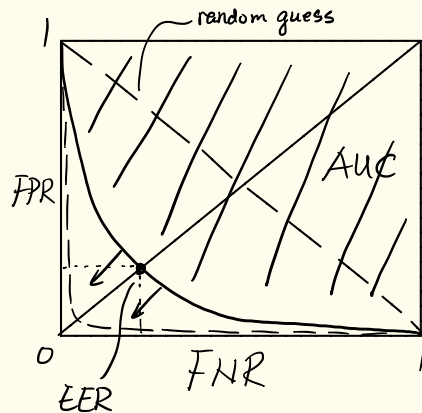
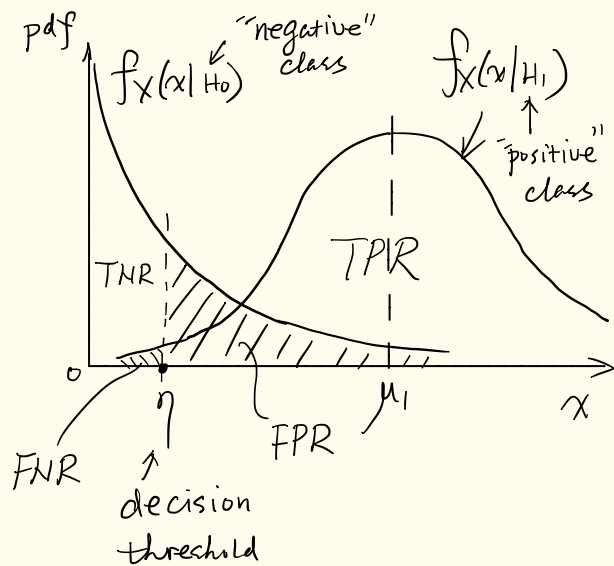


$$\text{Decision rule/classifier } \hat{C}(x) = \begin{cases} 1, & x \geq \eta, \\ 0, & x < \eta. \end{cases}$$

$$\begin{aligned}
 \text{FPR}(\eta) &= \int_{\eta}^{\infty} f_X(x|H_0) dx \\
 &= 1 - F_0(\eta) \\
 &= 1 - (1 - e^{-\lambda\eta}) = e^{-\lambda\eta}
 \end{aligned}$$

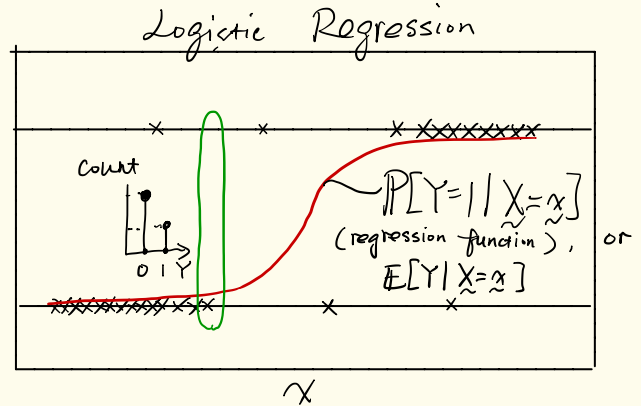
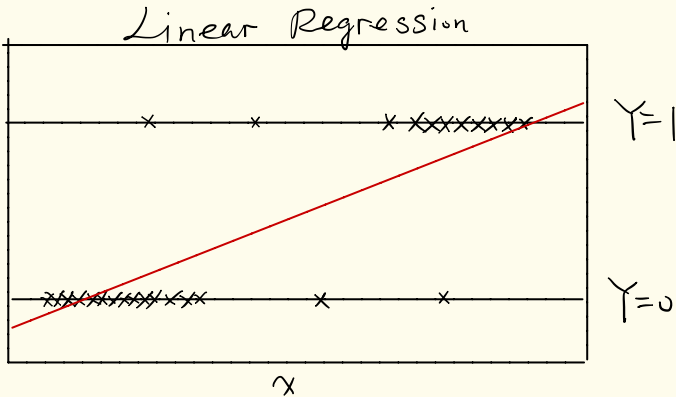
$$\begin{aligned}
 \text{FNR}(\eta) &= \int_{-\infty}^{\eta} f_X(x|H_1) dx \\
 &= F_1(\eta) = \Phi\left(\frac{x - \mu_1}{\sigma_1}\right)
 \end{aligned}$$

- ROC Curve $\stackrel{\text{def}}{=} \{(\text{FNR}(\eta), \text{FPR}(\eta))\}_{\eta}$
- EER: Equal Error Rate $\stackrel{\text{def}}{=} \text{FNR}(\eta) = \text{FPR}(\eta)$
- AUC: Area Under Curve $\stackrel{\text{def}}{=} \int_0^1 \text{TPR}(\text{FNR}) d(\text{FNR})$



Logistic Regression (LR) ESL 4.4, Murphy 8.1-3

Why logistic for binary data?



Need a nonlinear mapping $g: [0, 1] \rightarrow \mathbb{R}$

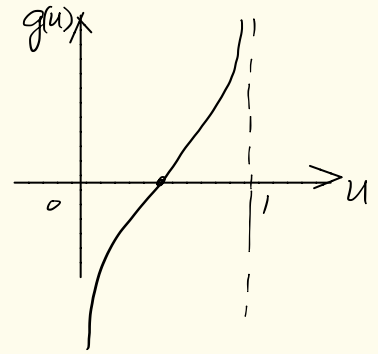
$P[Y=1 | X=\tilde{x}_i] \mapsto \tilde{x}_i^T \beta$

$\stackrel{\text{def}}{=} p(\tilde{x}_i)$

Not pdf/pmf, Not likelihood of \tilde{x}_i

— "logit": $g(u) = \log\left(\frac{u}{1-u}\right)$

— "Probit": $g(u) = \Phi^{-1}$ (historically)
inverse of CDF for Gaussian



Model:

odds log odds

$$\bullet \log\left(\frac{\overbrace{P[Y=1|x]}^{\text{odds}}}{\overbrace{P[Y=0|x]}^{\text{odds}}}\right) = \log\left(\frac{\overbrace{p(x)}^{\text{log odds}}}{1-p(x)}\right) = \beta_0 + \beta_1 x$$

$$\bullet g(\underbrace{p(x)}_u) = \beta_0 + \beta_1 x, \text{ where } g(u) = \log\left(\frac{u}{1-u}\right), \text{ or}$$

$$\bullet p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

Parameter Estimation for Logistic Regression

To estimate (β_0, β_1) . Training dataset $\{(x_i, y_i)\}_{i=1}^n$

How? Use Bayes classifier $\max p(y|x)$,
↑ predictor ↑ label.

which becomes MLE for (β_0, β_1) , treating labels $\{y_i\}$

as the "data" for MLE purposes: $\max_{\beta} \prod_{i=1}^n P[Y=y_i | X=x_i]$

$$L(\beta) = \prod_{i=1}^n P[Y=y_i | X=x_i]$$

$$= \prod_{i=1}^n p(x_i)^{y_i} [1-p(x_i)]^{1-y_i} \quad \beta = (\beta_0, \beta_1)$$

$$\ln L(\beta) = \ell(\beta) = \sum_{i=1}^n y_i \log p(x_i) + (1-y_i) \log (1-p(x_i)) \quad \leftarrow \text{"cross entropy"}$$

$$= \sum \log (1-p(x_i)) + y_i \log \frac{p(x_i)}{1-p(x_i)} \quad p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$\beta_0 + \beta_1 x$

$$= \sum -\ln(1 + e^{\beta_0 + \beta_1 x_i}) + \sum y_i (\beta_0 + \beta_1 x_i)$$

$$\frac{\partial l(\beta)}{\partial \beta_1} = -\sum \frac{1}{1 + e^{\beta_0 + \beta_1 x_i}} \cdot e^{\beta_0 + \beta_1 x_i} \cdot x_i + \sum y_i x_i$$

$$= -\sum p(x_i) x_i + \sum y_i x_i = \sum [y_i - \underbrace{p(x_i; \beta_0, \beta_1)}_{\text{nonlinear, needs IRLS}}] x_i$$

$$= \left. \begin{array}{l} 0 \\ \beta = \hat{\beta} \end{array} \right\} \text{IRLS}$$

Btw, for linear regression, using MLE (= LS when Gaussian)

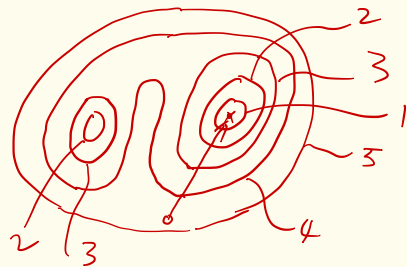
$$\min_{\beta_1} \sum_{i=1}^n (y_i - \beta_1 x_i)^2$$

$$\frac{\partial J}{\partial \beta_1} = \sum 2(y_i - \beta_1 x_i) \cdot (-x_i) = \left. \begin{array}{l} 0 \\ \beta_1 = \hat{\beta}_1 \end{array} \right\} \Rightarrow \sum [y_i - \hat{\beta}_1 x_i] x_i = 0$$

Optimization 101:

Goal $\min_{\tilde{x}} f(\tilde{x})$

f not necessarily convex



Approaches: 1. gradient / steepest descent (1st-order)

2. Newton's method (2nd-order)

3. Quasi-Newton (1st-order) "BFGS"

Algorithm: start w/ $x^{(0)}$

for $k=0, 1, 2, \dots$

compute $g^{(k)} = \nabla f(x^{(k)}) \leftarrow$ gradient (vector)

if steepest == 1

$$d^{(k)} = -g^{(k)}$$

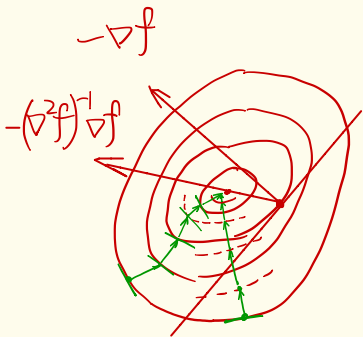
elseif Newton == 1

$$H^{(k)} = \nabla^2 f(x^{(k)}) \leftarrow \text{Hessian (matrix)}$$

$$d^{(k)} = -H^{(k)-1} g^{(k)} \quad \text{step size}$$

$$\text{Update: } x^{(k+1)} = x^{(k)} + \alpha_k \cdot d^{(k)}$$

Stop when $\|\nabla f(x^{(k)})\| < \epsilon$



$$g = \nabla_{\beta} l(\beta) = \sum_{i=1}^n (y_i - p(\tilde{x}_i; \beta)) \tilde{x}_i = X^T (y - p)$$

$$H = \nabla_{\beta}^2 l(\beta) = - \sum_{i=1}^n p(\tilde{x}_i) [1 - p(\tilde{x}_i)] \tilde{x}_i \tilde{x}_i^T = -X^T W X$$

where $W = \text{diag}(p(\tilde{x}_i)[1 - p(\tilde{x}_i)], i=1, \dots, n)$



$$\textcircled{1} \quad \beta^{\text{new}} \leftarrow \beta^{\text{old}} - H^{-1} g$$

$$= \beta^{\text{old}} + (X^T W X)^{-1} X^T (y - p)$$

$$= (X^T W X)^{-1} X^T W \left[X \beta^{\text{old}} + W^{-1} (y - p) \right]$$

$$= (X^T W X)^{-1} X^T W \tilde{z} \stackrel{\text{def}}{=} \tilde{z} \quad \text{"adjusted response"}$$

$$\Leftrightarrow \beta^{\text{new}} = \underset{\beta}{\text{argmin}} (\tilde{z} - X\beta)^T W (\tilde{z} - X\beta) \quad \text{Weighted LS (WLS).}$$

$\textcircled{2}$ Update W , or reweight

Prove yourself: $\hat{\beta}_{\text{WLS}} = (X^T W X)^{-1} X^T W \tilde{z}$

Iteratively
Reweighted
LS (IRLS)

Multinomial / Multiclass Logistic Regression

$$\log \frac{\mathbb{P}[Y=1|\underline{x}]}{\mathbb{P}[Y=k|\underline{x}]} = \beta_1^T \underline{x} \quad \left. \vphantom{\log \frac{\mathbb{P}[Y=1|\underline{x}]}{\mathbb{P}[Y=k|\underline{x}]}} \right\} (k-1) \text{ log odds, } Y \in \{1, \dots, k\}$$

$$\vdots$$

$$\log \frac{\mathbb{P}[Y=k-1|\underline{x}]}{\mathbb{P}[Y=k|\underline{x}]} = \beta_{k-1}^T \underline{x}$$

$$\mathbb{P}[Y=k|X=\underline{x}] = \frac{\exp(\beta_k^T \underline{x})}{1 + \sum_{l=1}^{k-1} \exp(\beta_l^T \underline{x})}, \quad k=1, \dots, k-1.$$

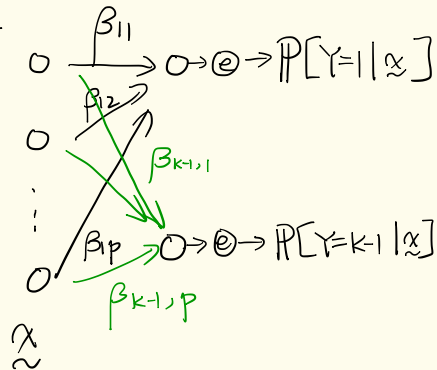
Soft(arg) max function: $\sigma_i(\underline{z}) = \frac{e^{\beta z_i}}{\sum_{j=1}^k e^{\beta z_j}}$

smooth version of "arg max".

Ex: $k=2$. $\sigma_1(\underline{z}) = \frac{e^{\beta z_1}}{e^{\beta z_1} + e^{\beta z_2}} = \frac{1}{1 + e^{\beta(z_2 - z_1)}}$

When β very large

$z_2 > z_1$ leads to $\begin{cases} \sigma_1 = 0 \\ \sigma_2 = 1 \end{cases}$



Generalized Linear Model (GLM)

For linear model , $g(\mu_i) = \beta_0 + \beta_1 x_i$, $g(\cdot)$ is identity

$$\mu_i = E[Y | x_i]$$

For logistic reg , $g(\mu_i) = \beta_0 + \beta_1 x_i$, $g(u) = \log \frac{u}{1-u}$

For Poisson dist , $g(\mu_i) = \beta_0 + \beta_1 x_i$, $g(u) = \log(u)$

$$Y_i \in \{0, 1, 2, \dots\}$$

↑
link func.

Exponential Family Distributions:

$$f(y; \theta) = \exp \left[\frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi) \right] , \text{ where } \phi \text{ is fixed}$$

↑
natural param

↑
scale param

Bernoulli:
$$P[Y=y] = p^y (1-p)^{1-y} = \exp[y \ln p + (1-y) \ln(1-p)]$$

$$= \exp \left[y \underbrace{\ln \frac{p}{1-p}}_{\theta} + \ln(1-p) \right]$$

Poisson:
$$P[Y=y] = \frac{\mu^y e^{-\mu}}{y!} = \exp \left[y \underbrace{\ln \mu}_{\theta} - \mu - \ln(y!) \right]$$

Gaussian:
$$f_Y(y) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(y-\mu)^2/2\sigma^2} \quad (\sigma^2 \text{ fixed, known})$$

$$= \exp \left[-\ln(\sqrt{2\pi}\sigma) - \frac{(y-\mu)^2}{2\sigma^2} \right]$$

$$= \exp \left[\frac{y \underbrace{\mu}_{\theta}}{\sigma^2} + C(y, \phi) \right]$$

$\underbrace{\hspace{10em}}_{\text{scale param}}$

$g(\cdot)$ such selected are called "canonical link function".

Exp family dist can use the same method, i.e., IRLS to find the best params:

$$\begin{aligned}\log \mathcal{L}(\beta) &= \log \left(\prod_{i=1}^n f(y_i; \beta) \right) \\ &= \sum_{i=1}^n \left[\frac{y_i \theta_i - b(\theta_i)}{a(\phi)} + c(y_i, \phi) \right]\end{aligned}$$

$$\frac{\partial}{\partial \beta_i} \log \mathcal{L}(\beta) = 0 \quad \text{for all } i.$$

Details about IRLS for GLM, see McCulloch 5.4.