# Naive Bayes X Logistic Regression

Let's say we want to classify a person into Male/Female based on hair length

Probabilistic Classifier →

$$P(C = C_K | x)$$

$$\Downarrow$$

$P(C = Male/Female\ given\ hair\ length)$

$C/y$ = Class Label , $x$ = feature vector

We need to find this probability.

If $P(C = M | x) > P(C = F | x)$; then we will
output Male

# Bayes Theorem $\rightarrow$

posterior      liklihood      prior knowledge

$$P(y_i \mid x_i) = \frac{P(x_i \mid y_i) \cdot P(y_i)}{P(x_i)}$$
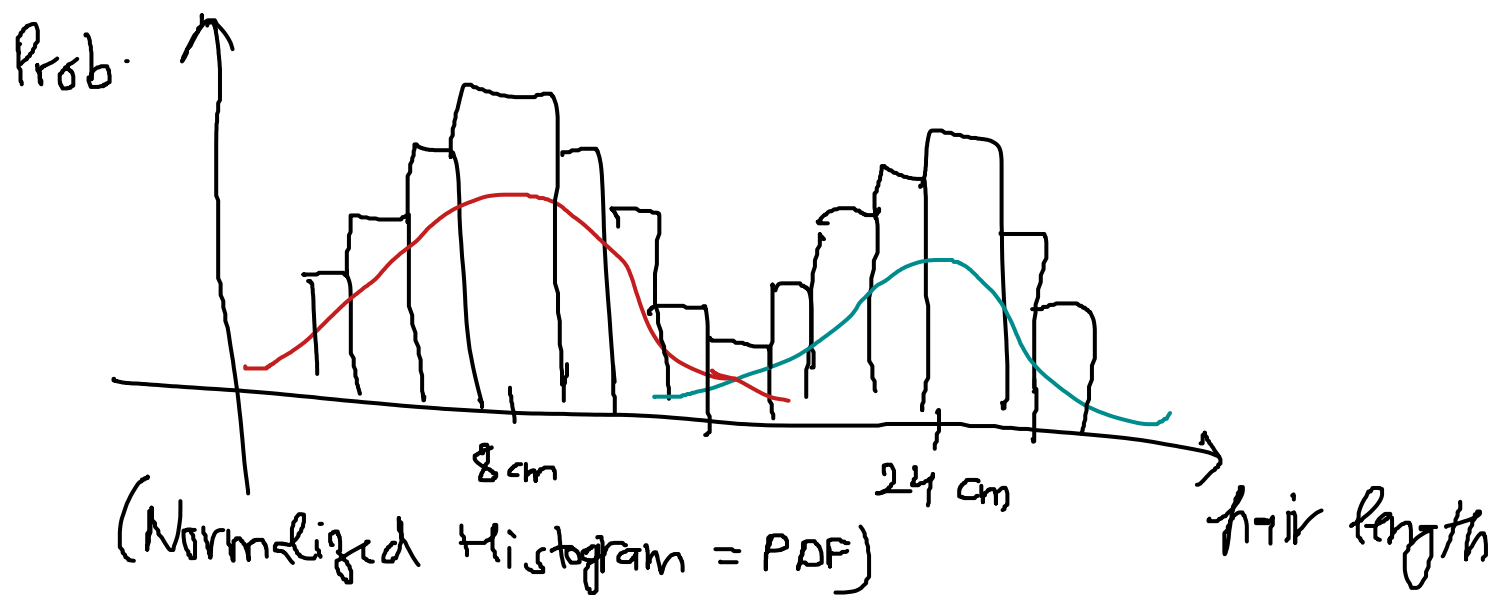
marginalization

Same for both classes
so we can ignore it

We know this
$$\frac{Male}{Total} , \frac{Female}{Total}$$

# Fitting a gaussian $\rightarrow$

$$p(x_i \mid \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \quad ; \{x_i\}_{i=1:N}$$



(Normalized Histogram = PDF)

To find the best gaussian we need to define a cost function.

Can you think of one?

# MLE →

We are maximising the prob. of N samples, while maximising the likelihood of a curve.

$(\theta = \{\mu, \sigma\})$

$$p(X|\theta) = p(x_1 \cdot x_2 \cdot x_3 \cdots x_N | \theta)$$

✱ iid assumptions → (independently & identically dist.)

$$p(X|\theta) = \prod_{i=1}^{N} p(x_i|\theta)$$

$$L(\theta) = \ln\left(p(X|\theta)\right) = \sum_{i=1}^{N} \ln\left(p(x_i|\theta)\right) \quad \left(\begin{array}{c} \log \text{ is} \\ \text{monotonic} \end{array}\right)$$

$$\hat{\theta} = \underbrace{\text{argmax}\left(L(\theta)\right)}_{\text{find optimal } \theta}$$

$$\nabla_\theta L = \sum_{i=1}^{N} \ln\left(p(x_i|\theta)\right) = 0$$

# Closed form solution →

$$p\left(x_i \mid \mu, \sigma\right) = \frac{1}{\sigma \sqrt{2\pi}} \, e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

$$\underbrace{\ln\left(p(\ )\right)}_{\text{loss function } L(\mu, \sigma)} = -N \ln\left(\sigma \sqrt{2\pi}\right) - \sum_{i=1}^{N} \frac{(x_i - \mu)^2}{2\sigma^2}$$

① $\dfrac{\partial L}{\partial \mu} = 0$ ; $\displaystyle\sum_{i=1}^{r} \frac{(x_i - \mu)}{\sigma^2} = 0$

$$\sum x_i - N\cdot\mu = 0$$

$$\boxed{\mu = \frac{\sum x_i}{N}} \quad (\text{mean})$$

② $\dfrac{\partial L}{\partial \sigma} = 0$ ; $\dfrac{-N \cdot \sqrt{2\pi}}{\sigma \sqrt{2\pi}} + \dfrac{\sum (x_i - \mu)^2}{\sigma^3}$

$$\boxed{\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{N}}} \quad \left(\begin{array}{l}\text{std} \\ \text{dev}\end{array}\right)$$

# Naive Bayes →

We can extend the above model to multiple features,

Given: 2D feature vector (hair length, voice pitch)

So, we will fit 4 gaussians

<span style="color:red">Male hair, Female hair, Male Voice, Female Voice</span>

Note that we still have 2 classes only, but our $x$ is now a matrix instead of vector.

$$P(C=C_k \mid x) = \frac{p(x_1 x_2 \mid C=C_k) \cdot P(C=C_k)}{p(x_1, x_2)}$$

$$= \frac{p(x_1 \mid C=C_k) \cdot p(x_2 \mid C=C_k) \cdot P(C=C_k)}{p(x_1, x_2)}$$

<span style="color:red">Naive Assumption → All features are independent</span>

<span style="color:blue">(Without Naive, we will have to fit multivariate gaussian)</span>

So, for M features →

$$X = {}_N\{\begin{bmatrix} \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}}$$

$\underbrace{\phantom{xxxxx}}_{M}$  $N \times M$

We can fit $2 \cdot M$ gaussians assuming all M features are mutually independent

Naive Bayes belongs to a class of models
known as Generative models.

Once we know the parameters of optimal gaussian
we can generate synthetic data points



We also have discriminative models which
cannot be used for generating data, because
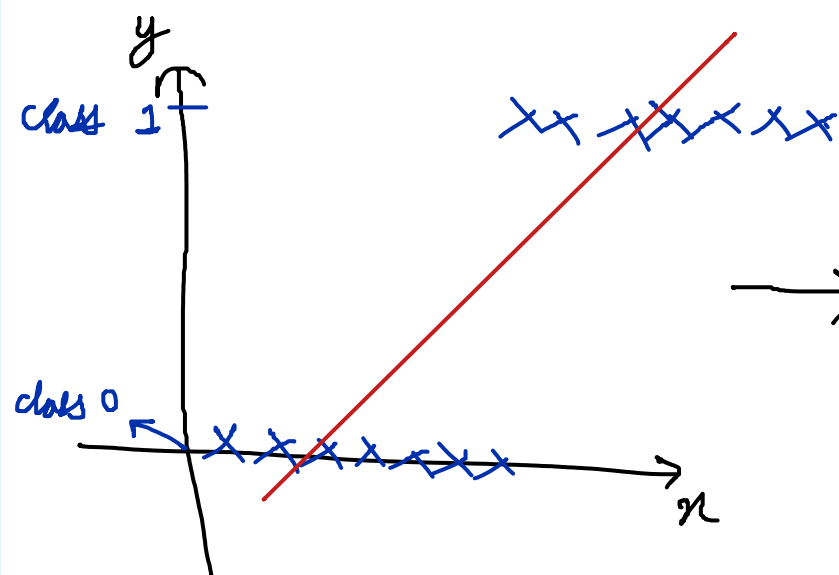they do not store information about distribution
of data.

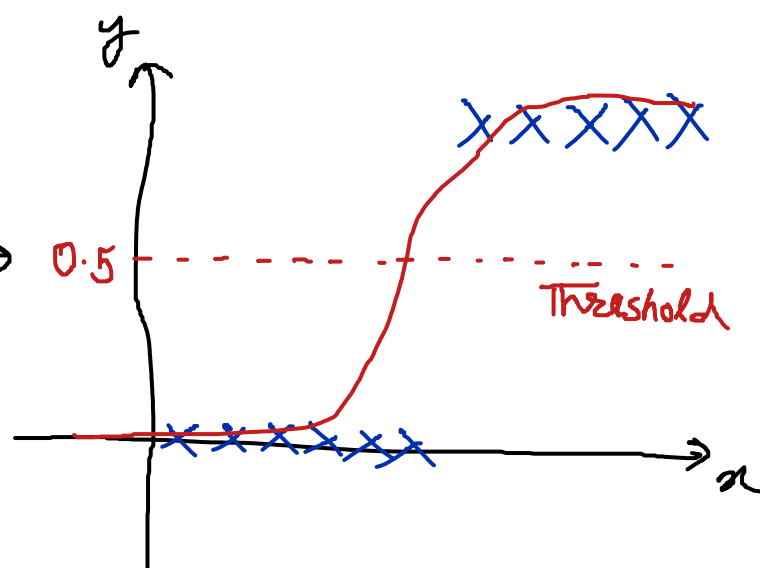What if our data is not gaussian?

# Logistic Regression →

Given $\vec{x}$, LR does the following →

① fit a linear classifier $(W^T x)$ similar to Linear Regression

② Apply Sigmoid function to output of ①

$$z = W^T x$$
$$(w_1 x_1 + w_b)$$

$$y = \sigma(z) = \frac{1}{1 + e^{-W^T x}}$$

Sigmoid converts $(-\infty, \infty) \to (0, 1)$ probabilities

# MLE (Bernoulli) →

Unlike Naive Bayes, here we will directly estimate the posterior probability.

$$\{M, F, M, M, M, F, F_{- - -}\} : N \text{ observations}$$

$$p = P(\text{Female} / 1)$$
$$1 - p = P(\text{Male} / 0)$$

$$P(Y ; p) = \prod_{i=1}^{N} P(y_i ; p) \qquad \left(\begin{array}{l}\text{assuming independant} \\ \text{bernoulli trials}\end{array}\right)$$

$$L(p) = \ln(P(y ; p)) = \sum \ln(P(y_i ; p))$$

From Bernoulli $= P(y_i ; p) = p^{y_i} \cdot (1-p)^{(1-y_i)}$

$$L(p) = \sum_{i=1}^{N} y_i \ln(p) + (1-y_i) \ln(1-p)$$

For our case, $p = \sigma(\alpha) = \dfrac{1}{1 + e^{-w^T x}}$

## Some Basic Results →

$$P(y=1 \mid x; w) = \frac{1}{1+e^{-w^T x}} = \frac{e^{w^T x}}{1+e^{w^T x}}$$

$$P(y=0 \mid x; w) = \frac{e^{-w^T x}}{1+e^{-w^T x}} = \frac{1}{1+e^{w^T x}}$$

Back to our loss function →

$$-\ell(w) \Rightarrow -\sum \left( y_i \ln\left(P(y=1 \mid x_i)\right) + (1-y_i)\ln\left(1-P(y_i=1 \mid x)\right)\right)$$

$$\Rightarrow -\sum \left( y_i \ln\left(\frac{1}{1+e^{-w^T x}}\right) + (1-y_i)\ln\left(\frac{1}{1+e^{w^T x}}\right)\right)$$

$$\Rightarrow \sum \left( y_i \ln\left(1+e^{-w^T x}\right) + (1-y_i)\ln\left(1+e^{w^T x}\right)\right)$$

$$\frac{\partial \ell}{\partial w} = \sum \left( y_i \cdot \frac{-x \cdot e^{-w^T x}}{1+e^{-w^T x}} + (1-y_i)\frac{x \cdot e^{w^T x}}{1+e^{w^T x}}\right)$$

$$= \sum \left( y_i \cdot \frac{-x \cdot e^{-w^T x}}{1+e^{-w^T x}} + (1-y_i)\frac{x}{1+e^{-w^T x}}\right)$$

$$= \sum \left( \frac{x_i - y_i x_i - x_i y_i e^{-W^T x_i}}{1 + e^{-W^T x}} \right)$$

$$= \sum \left( \frac{x_i - y_i x_i \left(1 + e^{-W^T x_i}\right)}{1 + e^{-W^T x_i}} \right)$$

$$= \sum \left( x_i \left( \sigma(x_i) - y_i \right) \right)$$

In matrix form $\Rightarrow$ $\boxed{\nabla_W L = X^T \left( \sigma(X) - \vec{y} \right)}$

---

Apply gradient descent to get optimal $w$.

$$X = \overset{x_1 \, x_2 \, x_3 \, 1}{\underset{N \times 4}{\begin{bmatrix} - & - & - & - \\ - & - & - & - \\ - & - & - & - \end{bmatrix}}} ; \quad \vec{y} = \underset{N \times 1}{\begin{bmatrix} \, \\ \, \end{bmatrix}} \quad \vec{w} = \underset{4 \times 1}{\begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_o \end{bmatrix}}$$

$$w \longrightarrow \text{random}$$

$$\underbrace{w}_{4 \times 1} \longrightarrow \underbrace{w}_{4 \times 1} - \alpha \cdot \underbrace{X^T \left( \underbrace{\sigma(x)}_{N \times 1} - \underbrace{y}_{N \times 1} \right)}_{4 \times 1}$$

$$\sigma(W^T X)$$
$$\approx X W$$
$$N \times 4 \cdot 4 \times 1$$
$$= N \times 1$$

Note $\rightarrow$ L.R. can also be derived from

Cross Entropy Loss,

$$\mathcal{L}_{CE} = -\sum y_i \log(p_i)$$

This is equivalent to Log loss or Logistic loss we derived earlier.

What if there are more than 2 classes?

$$\text{Softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}}$$

Converts vector of K numbers into PDF of K outcomes